

A

03/15/99
JC557 U.S. PTO

DOCKET NO. : MSFT-0038/36765.2

PATENT

JC525 U.S. PTO
09/268146
03/15/99

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In Re Application of:

David B. Lomet; Mark R. Tuttle

Serial No.: Not yet assigned

Group Art Unit: Not yet assigned

Filing Date: Herewith

Examiner: Not yet assigned

For: DATABASE COMPUTER SYSTEM USING LOGICAL LOGGING TO
EXTEND RECOVERY

EXPRESS MAIL LABEL NO: EL140211439US

DATE OF DEPOSIT: March 15, 1999

Box ☒ Patent Application
☐ Provisional ☐ Design ☐ Sequence

Assistant Commissioner for Patents
Washington DC 20231

Sir:

PATENT APPLICATION TRANSMITTAL LETTER

Transmitted herewith for filing, please find

☒ A Utility Patent Application under 37 C.F.R. 1.53(b).

It is a continuing application, as follows:

☐ continuation ☐ divisional ☒ continuation-in-part of prior application number
08/832,870, filed April 4, 1997, which is a continuation-in-part of application number _
08/814,808, filed March 10, 1997 or application number 08/813,982, filed March 10,
1997, which has issued as U.S. Patent No. 5,870,763.

- ☐ A Provisional Patent Application under 37 C.F.R. 1.53(c).
- ☐ A Design Patent Application (submitted in duplicate).

Including the following:

- ☐ Provisional Application Cover Sheet.
- ☒ New or Revised Specification, including pages 1 to 90 containing:

- ☒ Specification
- ☒ Claims
- ☒ Abstract
- ☐ Substitute Specification, including Claims and Abstract.

☐ The present application is a continuation application of Application No. _____ filed _____. The present application includes the Specification of the parent application which has been revised in accordance with the amendments filed in the parent application. Since none of those amendments incorporate new matter into the parent application, the present revised Specification also does not include new matter.

☐ The present application is a continuation application of Application No. _____ filed _____, which in turn is a continuation-in-part of Application No. _____ filed _____. The present application includes the Specification of the parent application which has been revised in accordance with the amendments filed in the parent application. Although the amendments in the parent C-I-P application may have incorporated new matter, since those are the only revisions included in the present application, the present application includes no new matter in relation to the parent application.

☐ A copy of earlier application Serial No. _____ Filed _____,

including Specification, Claims and Abstract (pages 1 - @@), to which no new matter has been added TOGETHER WITH a copy of the executed oath or declaration for such earlier application and all drawings and appendices. Such earlier application is hereby incorporated into the present application by reference.

- ☐ Please enter the following amendment to the Specification under the Cross-Reference to Related Applications section (or create such a section) : "This Application is a ☐ continuation or ☐ divisional of Application Serial No. _____ filed _____."
- ☐ Signed Statement attached deleting inventor(s) named in the prior application.
- ☐ A Preliminary Amendment.
- ☒ Twenty-Five (25) Sheets of ☒ Formal ☐ Informal Drawings.
- ☐ Petition to Accept Photographic Drawings.
- ☐ Petition Fee
- ☒ An ☐ Executed ☒ Unexecuted Declaration or Oath and Power of Attorney.
- ☐ An Associate Power of Attorney.
- ☐ An ☐ Executed ☐ Copy of Executed Assignment of the Invention to _____
- ☐ A Recordation Form Cover Sheet.
- ☐ Recordation Fee - \$40.00.
- ☐ The prior application is assigned of record to _____
- ☐ Priority is claimed under 35 U.S.C. § 119 of Patent Application No. _____ filed _____

_____ in _____ (country).

☐ A Certified Copy of each of the above applications for which priority is claimed:

☐ is enclosed.

☐ has been filed in prior application Serial No. _____ filed _____ .

☐ An ☐ Executed or ☐ Copy of Earlier Statement Claiming Small Entity Status under 37 C.F.R. 1.9 and 1.27

☐ is enclosed.

☐ has been filed in prior application Serial No. _____ filed _____ , said status is still proper and desired in present case.

☐ Diskette Containing DNA/Amino Acid Sequence Information.

☐ Statement to Support Submission of DNA/Amino Acid Sequence Information.

☐ The computer readable form in this application _____, is identical with that filed in Application Serial Number _____, filed _____. In accordance with 37 CFR 1.821(e), please use the ☐ first-filed, ☐ last-filed or ☐ only computer readable form filed in that application as the computer readable form for the instant application. It is understood that the Patent and Trademark Office will make the necessary change in application number and filing date for the computer readable form that will be used for the instant application. A paper copy of the Sequence Listing is ☐ included in the originally-filed specification of the instant application, ☐ included in a separately filed preliminary amendment for incorporation into the specification.

☒ Information Disclosure Statement.

☒ Attached Form 1449.

☒ Copies of each of the references listed on the attached Form PTO-1449 are enclosed herewith.

☐ A copy of Petition for Extension of Time as filed in the prior case.

☐ Appended Material as follows: _____.

☒ Return Receipt Postcard (should be specifically itemized).

☐ Other as follows: _____

 _____.

FEE CALCULATION:

☐ Cancel in this application original claims _____ of the prior application before calculating the filing fee. (At least one original independent claim must be retained for filing purposes.)

				SMALL ENTITY		NOT SMALL ENTITY	
				RATE	FEE	RATE	FEE
PROVISIONAL APPLICATION				\$75.00	\$	\$150.00	\$
DESIGN APPLICATION				\$155.00	\$	\$310.00	\$
UTILITY APPLICATIONS BASE FEE				\$380.00	\$	\$760.00	\$760.00
UTILITY APPLICATION; ALL CLAIMS CALCULATED AFTER ENTRY OF ALL AMENDMENTS							
	No. Filed	No. Extra					
TOTAL CLAIMS	40 - 20 =	20		\$9 each	\$	\$18 each	\$360.00
INDEP. CLAIMS	8 - 3 =	5		\$39 each	\$	\$78 each	\$390.00
FIRST PRESENTATION OF MULTIPLE DEPENDENT CLAIM				\$130	\$	\$260	\$ 0
ADDITIONAL FILING FEE					\$		\$ 0
TOTAL FILING FEE DUE					\$		\$1,510.00

☒ A Check is enclosed in the amount of \$ 1,510.00.

☒ The Commissioner is authorized to charge payment of the following fees and to refund

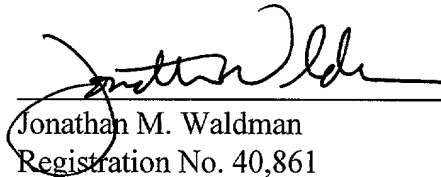
any overpayment associated with this communication or during the pendency of this application to deposit account 23-3050. This sheet is provided in duplicate.

- ☐ The foregoing amount due.
- ☒ Any additional filing fees required, including fees for the presentation of extra claims under 37 C.F.R. 1.16.
- ☒ Any additional patent application processing fees under 37 C.F.R. 1.17 or 1.20(d).
- ☐ The issue fee set in 37 C.F.R. 1.18 at the mailing of the Notice of Allowance.

- ☒ The Commissioner is hereby requested to grant an extension of time for the appropriate length of time, should one be necessary, in connection with this filing or any future filing submitted to the U.S. Patent and Trademark Office in the above-identified application during the pendency of this application. The Commissioner is further authorized to charge any fees related to any such extension of time to deposit account 23-3050. This sheet is provided in duplicate.

SHOULD ANY DEFICIENCIES APPEAR with respect to this application, including deficiencies in payment of fees, missing parts of the application or otherwise, the United States Patent and Trademark Office is respectfully requested to promptly notify the undersigned.

Date: **March 15, 1999**


Jonathan M. Waldman
Registration No. 40,861

Woodcock Washburn Kurtz
Mackiewicz & Norris LLP
One Liberty Place - 46th Floor
Philadelphia PA 19103
Telephone: (215) 568-3100
Facsimile: (215) 568-3439

© 1997 WWKMN

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Database Computer System Using Logical Logging
To Extend Recovery**

Inventors:

David B. Lomet

Mark R. Tuttle

RELATED APPLICATIONS

This is a continuation-in-part of U.S. Patent Application Serial Number 08/832,870, which was filed April 4, 1997, which is a continuation-in-part of U.S. Patent Application Serial Number 08/814,808, or U.S. Patent Application Serial Number 08/813,982, which has issued as U.S. Patent Number 5,870,763, which were both filed March 10, 1997 in the name of David B. Lomet, and are both assigned to Microsoft Corporation.

TECHNICAL FIELD

This invention relates to database computer systems and applications that execute on them. More particularly, this invention relates to methods for recovering from system crashes in a manner that ensures that the applications themselves persist across the crash.

BACKGROUND OF THE INVENTION

Computer systems occasionally crash. A "system crash" is an event in which the computer quits operating the way it is supposed to operate. Common causes of system crashes include power outage, application operating error, and computer goblins (i.e., unknown and often unexplained malfunctions that tend to plague even the best-devised systems and applications). System crashes are unpredictable, and hence, essentially impossible to anticipate and prevent.

A system crash is at the very least annoying, and may result in serious or irreparable damage. For standalone computers or client workstations, a local system crash typically results in loss of work product since the last save interval. The user is inconvenienced by having to reboot the computer and redo the lost work. For servers and

larger computer systems, a system crash can have a devastating impact on many users, including both company employees as well as its customers.

Being unable to prevent system crashes, computer system designers attempt to limit the effect of system crashes. The field of study concerning how computers recover from system crashes is known as "recovery." Recovery from system crashes has been the subject of much research and development.

In general, the goal of redo recovery is to return the computer system after a crash to a previous and presumed correct state in which the computer system was operating immediately prior to the crash. Then, transactions whose continuations are impossible can be aborted. Much of the recovery research focuses on database recovery for database computer systems, such as network database servers or mainframe database systems. Imagine the problems caused when a large database system having many clients crashes in the midst of many simultaneous operations involving the retrieval, update, and storage of data records. Database system designers attempt to design the database recovery techniques which minimize the amount of data lost in a system crash, minimize the amount of work needed following the crash to recover to the pre-crash operating state, and minimize the performance impact of recovery on the database system during normal operation.

Fig. 1 shows a database computer system 20 having a computing unit 22 with processing and computational capabilities 24 and a volatile main memory 26. The volatile main memory 26 is not persistent across crashes and hence is presumed to lose all of its data in the event of a crash. The computer system also has a non-volatile or stable database 28 and a stable log 30, both of which are contained on stable memory devices, e.g. magnetic disks, tapes, etc., connected to the computing unit 22. The stable database 28 and log 30 are presumed to persist across a system crash. The persistent database 28

1 and log 30 can be combined in the same storage, although they are illustrated separately
2 for discussion purposes.

3 The volatile memory 26 stores one or more applications 32, which execute on the
4 processor 24, and a resource manager 34. The resource manager 34 includes a volatile
5 cache 36, which temporarily stores data destined for the stable database 28. The data is
6 typically stored in the stable database and volatile cache in individual units, such as
7 "pages." A cache manager 38 executes on the processor 24 to manage movement of data
8 pages between the volatile cache 36 and the stable database 28. In particular, the cache
9 manager 38 is responsible for deciding which data pages should be moved to the stable
10 database 28 and when the data pages are moved. Data pages which are moved from the
11 cache to the stable database are said to be "flushed" to the stable state. In other words,
12 the cache manager 38 periodically flushes the cached state of a data page to the stable
13 database 28 to produce a stable state of that data page which persists in the event of a
14 crash, making recovery possible.

15 The resource manager 34 also has a volatile log 40 which temporarily stores
16 computing operations to be moved into the stable log 30. A log manager 42 executes on
17 the processor 24 to manage when the operations are moved from the volatile log 40 to the
18 stable log 30. The transfer of an operation from the volatile log to the stable log is known
19 as a log flush.

20 During normal operation, an application 32 executes on the processor 24. The
21 resource manager receives requests to perform operations on data from the application.
22 As a result, data pages are transferred to the volatile cache 36 on demand from the stable
23 database 28 for use by the application. During execution, the resource manager 34 reads,
24 processes, and writes data to and from the volatile cache 36 on behalf of the application.
25

1 The cache manager 38 determines, independently of the application, when the cached
2 Data State is flushed to the stable database 28.

3 Concurrently, the operations being performed by the resource manager on behalf
4 of the application are being recorded in the volatile log 40. The log manager 42
5 determines, as guided by the cache manager and the transactional requirements imposed
6 by the application, when the operations are posted as log records on the stable log 30. A
7 logged operation is said to be "installed" when the versions of the pages containing the
8 changes made by the operation have been flushed to the stable database.

9 When a crash occurs, the application state (i.e., address space) of any executing
10 application 32, the data pages in volatile cache 36, and the operations in volatile log 40
11 all vanish. The computer system 20 invokes a recovery manager which begins at the last
12 flushed state on the stable database 28 and replays the operations posted to the stable log
13 30 to restore the database of the computer system to the state as of the last logged
14 operation just prior to the crash.

15 Explaining how to recover from a system crash requires answering some
16 fundamental questions.

- 18 1. How can the designer be sure that recovery will succeed?
- 19 2. How can the stable state be explained in terms of what operations have been
20 installed and what operations have not?
- 21 3. How should recovery choose the operations to redo in order to recover an
22 explainable state?
- 23 4. How should the cache manager install operations via its flushing of database
24 pages to the stable state in order to keep the state explainable, and hence
25 recoverable?

1
2 The answers to these questions can be found in delicately balanced and highly
3 interdependent decisions that a system designer makes.

4 One prior art approach to database recovery is to require the cache manager to
5 flush the entire cache state periodically. The last such flushed state is identified in a
6 "checkpoint record" that is inserted into the stable log. During recovery, a redo test is
7 performed to determine whether a logged operation needs to be redone to help restore the
8 system to its pre-crash state. The redo test is simply whether an operation follows the last
9 checkpoint record on the log. If so (meaning that a later operation occurred and was
10 posted to the stable log, but the results of the operation were not installed in the stable
11 database), the computer system performs a redo operation using the log record.

12 This simple approach has a major drawback in that writing every change of the
13 cached state out to the stable database 28 is practically unfeasible. It involves a high
14 volume of input/output (I/O) activity that consumes a disproportionate amount of
15 processing resources and slows the system operation. It also requires atomic flushing of
16 multiple pages, which is a troublesome complication. This was the approach used in
17 System R., described in: Gray, McJones, et al, "The Recovery Manager of the System R
18 Database Manager," ACM Computing Surveys 13,2 (June, 1981) pages 223-242.

19 Another prior art approach to database recovery, which is more widely adopted
20 and used in present-day database systems, involves segmenting data from the stable
21 database into individual fixed units, such as pages. Individual pages are loaded into the
22 volatile cache and logged resource manager operations can read and write only within the
23 single pages, thereby modifying individual pages. The cache manager does not flush the
24 page after every incremental change.

1 Each page can be flushed atomically to the stable database, and independently of
2 any other page. Intelligently flushing a page after several updates have been made to the
3 page produces essentially the same result as flushing each page after every update is
4 made. That is, flushing a page necessarily includes all of the incremental changes made
5 to that page leading up to the point when the flushing occurs.

6 The cache manager assigns a monotonically increasing state ID to the page each
7 time the page is updated. During recovery, each page is treated as if it were a separate
8 database. Resource manager operations posted to the stable log are also assigned a state
9 ID. A redo test compares, for each page, the state ID of a stable log record with the state
10 ID of the stable page. If the log record state ID is greater than the state ID of the stable
11 page (meaning that one or more operations occurred later and were recorded in the stable
12 log, but the page containing updates caused by the later operations was not yet flushed to
13 the stable database), the computer system performs a redo operation using the last stable
14 page and the operations posted to the stable log that have state IDs higher than the state
15 ID of the stable page.

16 While these database recovery techniques are helpful for recovering data, in the
17 database, the recovery techniques offer no help to recovering an application from a
18 system crash. Usually all active applications using the database are wiped out during a
19 crash. Any state in an executing application is erased and cannot usually be continued
20 across a crash.

21 Fig. 2 shows a prior art system architecture of the database computer system 20.
22 The applications 32(1)-32(N) execute on the computer to perform various tasks and
23 functions. During execution, the applications interact with the resource manager 26, with
24 each other, and with external devices, as represented by an end user terminal 44. The
25 application states can change as a result of application execution, interaction with the

1 resource manager 26, interaction with each other, and interaction with the terminal 44. In
2 conventional systems, the application states of the executing applications 32(1)-32(N) are
3 not captured. There is no mechanism in place to track the application state as it changes,
4 and hence, there is no way to recover an application from a crash which occurs during its
5 execution.

6 When the application is simple and short, the fact that applications are not
7 recoverable is of little consequence. For example, in financial applications like
8 debit/credit, there may be nothing to recover that was not already captured by the state
9 change within the stable database. But this might not always be the case. Long running
10 applications, which frequently characterize workflow systems, present problems. Like
11 long transactions that are aborted, a crash interrupted application may need to be re-
12 scheduled manually to bring the application back online. Applications can span multiple
13 database transactions whereby following a system crash, the system state might contain
14 an incomplete execution of the application. Cleanly coping with partially completed
15 executions can be very difficult. One cannot simply re-execute the entire activity because
16 the partially completed prior execution has altered the state. Further, because some state
17 changes may have been installed in the stable database, one cannot simply undo the entire
18 activity because the transactions are guaranteed by the system to be persistent. The
19 transactions might not be undoable in any event because the system state may have
20 changed in an arbitrary way since they were executed.

21 Accordingly, there is a need for recovery procedures for preserving applications
22 across a system crash. Conceptually, the entire application state (i.e., the address space)
23 could be posted to the stable log after each operation. This would permit immediate
24 recovery of the application because the system would know exactly, from the last log
25 entry for the application, the entire application state just prior to crash. Unfortunately, the

1 address space is typically very large and continuously logging such large entries is too
2 expensive in terms of I/O processing resources and the large amounts of memory required
3 to hold successive images of the application state.

4 There are several prior art techniques that have been proposed for application
5 recovery. All have difficulties that restrict their usefulness. One approach is to make the
6 application "stateless." Between transactions, the application is in its initial state or a
7 state internally derived from the initial state without reference to the persistent state of the
8 database. If the application fails between transactions, there is nothing about the
9 application state that cannot be re-created based on the static state of the stored form of
10 the application. Should the transaction abort, the application is replayed, thereby re-
11 executing the transaction as if the transaction executed somewhat later. After the
12 transaction commits, the application returns to the initial state. This form of transaction
13 processing is described by Gray and Reuter in a book entitled, Transaction Processing:
14 Concepts and Techniques, Morgan Kaufmann (1993), San Mateo, CA.

15 Another approach is to reduce the application state to some manageable size and
16 use a recoverable resource manager to store it. The resource manager might be a database
17 or a recoverable queue. Reducing state size can be facilitated by the use of a scripting
18 language for the application. In this case, the script language interpreter stores the entire
19 application state at well-chosen times so that failures at inappropriate moments survive,
20 and the application execution can continue from the saved point.

21 Another technique is to use a persistent programming language that logs updates
22 to a persistent state. The idea is to support recoverable storage for processes. When the
23 entire state of the application is contained in recoverable storage, the application itself can
24 be recovered. Recoverable storage has been handled by supporting a virtual memory
25 abstraction with updates to memory locations logged during program execution. If the

1 entire application state is made recoverable, a very substantial amount of logging activity
2 arises. This technique is described in the following publications: Chang and Mergen,
3 "801 Storage: Architecture and Programming," ACM Trans. on Computer Systems, 6, 1
4 (Feb. 1988) pages 28-50; and Haskin et al., "Recovery Management in QuickSilver,"
5 ACM Trans. on Computer Systems, 6,1 (Feb. 1988) pages 82-108.

6 Another approach is to write persistent application checkpoints at every resource
7 manager interaction. The notion here is that application states in between resource
8 manager interactions can be re-created from the last such interaction forward. This is the
9 technique described by Bartlett, "A NonStop Kernel," Proc. ACM Symp. on Operating
10 System Principles (1981) pages 22-29 and Borg et al. "A Message System Supporting
11 Fault Tolerance," Proc. ACM Symp. on Operating System Principles (Oct. 1983) Bretton
12 Woods, NH pages 90-99. The drawback with this approach is that short code sequences
13 between interactions can mean frequent checkpointing of very large states as the state
14 changes are not captured via operations, although paging techniques can be used to
15 capture the differences between successive states at, perhaps, page level granularity.

16 The inventor has developed an improved recovery technique that breaks apart
17 flush dependencies that require atomic flushing of more than one object simultaneously.
18 This enables an ordered flushing sequence of first flushing a first object and then flushing
19 a second object, rather than having to flush both the first and second objects
20 simultaneously and atomically.

21 22 SUMMARY OF THE INVENTION

23 This invention concerns a database computer system and method for making
24 applications recoverable from system crashes. The application state (i.e., address space)
25 is treated as a single object that can be atomically flushed in a manner akin to flushing

individual pages in database recovery techniques. And like the pages of the database, log records describing application state changes are posted on the stable log before application state is flushed.

To enable this monolithic treatment of the application, executions performed by the application are mapped to loggable operations which are posted to the stable log. Any modifications to the application state are accumulated and the application state is flushed from time to time to stable storage using an atomic write procedure. Flushing the application state to stable storage effectively installs the application operations logged in the stable log. Since the application state can be very large, a procedure known as "shadowing" can be used to atomically flush the entire application state. As a result, the application recovery integrates with database recovery, and substantially reduces the need for checkpointing applications, i.e. logging or flushing the entire application state. According to one implementation, a database computer system has a processing unit, a volatile main memory that does not persist across a system crash, and a stable memory that persists across a system crash. The volatile memory includes a volatile cache which maintains cached states of the application address space and data records and a volatile log which tracks the operations performed by the computer system. The stable memory includes a stable database which stores stable states of the application address space and data records and a stable log which holds a stable version of the log records that describe state changes to the stable database.

The database computer system has at least one application which executes from the main memory on the processing unit. A resource manager is stored in main memory and mediates all interaction between the application and the external world (e.g., user terminal, data file, another application, etc.). During execution, the internal state changes of the application are not visible to the outside world. However, each time the

1 application interacts with the resource manager, either the application state is exposed or
2 the application senses the external state. The resource manager tags the application states
3 at these interaction points by assigning them state IDs. Application operations are
4 defined that produce the transitions between these application states. These operations
5 are immediately entered into the volatile log, and subsequently posted to the stable log.

6 The application state is treated as a single object that can be atomically flushed to
7 the stable database. In addition, the application operations often cause changes to the
8 data pages, records, or other types of objects stored in the volatile cache. The modified
9 objects that result from application operations are from time to time flushed to the stable
10 database. The flushed application states and objects are assigned state IDs to identify
11 their place in the execution sequence. Flushing the application object effectively installs
12 all the operations, updating the application operations that are in the stable log which
13 have earlier state IDs.

14 In the event of a system failure, the database computer system begins with the
15 stable database state and replays the stable log to redo certain logged application
16 operations. The database computer system redoes a logged application operation if its
17 state ID is later in series than the state ID of the most recently flushed or already partially
18 recovered application state.

19 Another aspect of this invention is to optimize the application read operation to
20 avoid writing the object data read to the log record. Posting the read values to the log is
21 helpful in one sense because the cache manager is not concerned about which sequence to
22 flush objects. Certain object states need not be preserved by a particular flushing order
23 because any data values obtained from an object which are needed to redo an application
24 operation are available directly from the stable log. However, posting objects to the log
25

1 often involves writing large amounts of data, and duplicating data found elsewhere on the
2 system.

3 The read optimizing technique eliminates posting the read values to the log by
4 substituting, for the read values, an identity of the location from where the values are read
5 and posting the identity instead of the values. However, the data is now only available
6 from the read object itself and hence, attention must be paid to the order in which objects
7 are flushed to stable storage. If objects are flushed out of proper sequence, a particular
8 state of an object may be irretrievably lost.

9 A cache manager has an object table which tracks the objects maintained in the
10 volatile cache. The object table includes fields to track dependencies among the objects.
11 In one implementation, the object table includes, for each object entry, a predecessor field
12 which lists all objects that must be flushed prior to the subject object, and a successor
13 field which lists all objects before which the subject object must be flushed. In another
14 implementation, the object table contains, for each object entry, a node field to store
15 dependencies in terms of their nodes in a write graph formulation.

16 Another aspect of this invention is to optimize the application write operation to
17 avoid posting large amounts of data to the log record. Posting the values to be written is
18 helpful in one sense because the cache manager is not concerned about which sequence to
19 flush objects. However, the process is inefficient and costly in terms of computational
20 resources.

21 The write optimization technique eliminates posting the write values to the log by
22 substituting, for those values, an identity of the object from where the values originate
23 and posting the identity instead of the values. While this reduces the amount of data to be
24 logged, the write optimization technique introduces dependencies between objects, and
25

often troubling “cycle” dependencies when the read optimization technique is also being used, which can require atomic and simultaneous flushing of multiple objects.

The cache manager tracks dependencies via the object table and is configured to recognize cycle dependencies. When a cycle dependency is realized, the cache manager initiates a blind write of one or more objects involved in the cycle to place the objects’ values on the stable log. This process breaks the cycle. Thereafter, the cache manager flushes the objects according to an acyclic flushing sequence that pays attention to any predecessor objects that first require flushing.

Therefore, in a database computer system having a non-volatile memory, a volatile main memory, and an application object which executes from the main memory, wherein the non-volatile memory includes a stable log, a computer-implemented method in accordance with the present invention comprises the following steps: executing the application object to perform operations which read data from, and write data to, a data object; posting to the stable log a log record for each operation involving the reading or writing of data, the log record containing a reference to either the application object or the data object to identify that referenced object as a source for the data that is read from or written to; establishing flush order dependencies between the application object and the data object, wherein some of the flush order dependencies become cyclic indicating a condition in which the application object should be flushed not later than the data object and the data object should be flushed not later than the application object; detecting a dependency cycle; and following detection of the dependency cycle, writing one of the application object or the data object to the stable log to break the dependency cycle to enable the application and data objects to be flushed sequentially according to an ordered flushing sequence. It should be noted that the technique of the present invention can be

used for breaking up atomic flush sets, regardless of how they arise (e.g., as a result of cyclic flush dependencies, as a result of an operation that updates two objects, etc.).

Preferably, the writing step comprises writing the data object to the stable log. More preferably, the method comprises the step of flushing the application object to the non-volatile memory after the data object is written to the stable log, and the method further comprises the step of flushing the data object to the non-volatile memory after the application object has been flushed to the non-volatile memory. The step of subsequently flushing the data object is to permit the object to be dropped from the cache. The value of the data object can be retrieved from its stable (non-volatile) storage location if it is needed.

In accordance with another aspect of the present invention, in a database computer system having a cache manager which occasionally flushes objects from a volatile main memory to a non-volatile memory to preserve those objects in the event of a system crash, and wherein a dependency cycle exists between at least two objects such that the two objects should be flushed simultaneously, a computer-implemented method comprises the following steps: detecting a dependency cycle; and writing one of the two objects to the stable log to break the dependency cycle to enable the two objects to be flushed to the non-volatile memory in a sequential manner according to an ordered flushing sequence. The method preferably comprises the step of flushing the objects according to the ordered flushing sequence after the writing step.

Thus, according to one aspect of the present invention, the acyclic flushing sequence is structured such that the object that is removed from the cycle dependency by the blind write is flushed to the stable database after the other object of the original cycle dependency. In other words, the object that is not removed from the cycle dependency by

1 the blind write is flushed to the stable database before the object that is removed from the
2 cycle dependency is flushed to the stable database.

3 Still another aspect of this invention is to optimize the recovery procedures
4 invoked following a system crash. During normal operation, each log record is assigned
5 a log sequence number (LSN). The cache manager maintains a recovery log sequence
6 number (rLSN) that identifies the first log record for an associated object at which to
7 begin replaying the operations during recovery. The cache manager occasionally flushes
8 an object to non-volatile memory to install the operations performed on the object. On
9 some occasions, the flushing of one object installs operations that wrote another data
10 object that has not yet been flushed (i.e., an object that is unexposed in the write graph,
11 meaning that its contents are not needed for recovery). The cache manager advances the
12 rLSN for both objects to identify subsequent log records that reflect the objects at states
13 in which the operations that previously wrote the states are installed in the non-volatile
14 memory.

15 During recovery, the recovery manager starts at the advanced rLSNs to avoid
16 replaying operations that are rendered unnecessary by subsequent operations, thereby
17 optimizing recovery.

18 19 **BRIEF DESCRIPTION OF THE DRAWINGS**

20 Fig. 1 is a diagrammatic illustration of a conventional database computer system.

21 Fig. 2 is a diagrammatic illustration of a system architecture of the conventional
22 database computer system.

23 Fig. 3 is a diagrammatic illustration of a database computer system according to
24 an implementation of this invention.
25

1 Fig. 4 is a diagrammatic illustration of a cache manager and non-volatile memory
2 used in the database computer system, and demonstrates aspects concerning atomic
3 installation of large application objects.

4 Fig. 5 is a diagrammatic illustration of a system architecture of the database
5 computer system that enables application recovery.

6 Fig. 6 is a diagrammatic illustration of application execution and interaction with
7 a resource manager in a manner which maps application execution to loggable logical
8 operations. Fig. 6 shows a logical execution operation.

9 Fig. 7 is a diagrammatic illustration similar to Fig. 6, but showing a logical read
10 operation.

11 Fig. 8 is a diagrammatic illustration similar to Fig. 6, but showing a logical write
12 operation.

13 Fig. 9 is a diagrammatic illustration showing a sequence of logical application
14 operations and how the operations are logged.

15 Fig. 10 is a diagrammatic illustration of the sequence of operations from Fig. 9,
16 which shows a read optimizing technique for logging operations and objects affected by
17 read operations.

18 Fig. 11 is a diagrammatic illustration of a cache manager with an object table for
19 tracking dependencies between data and application objects.

20 Fig. 12 is a write graph that illustrates a read-write dependency between an
21 application object and a data object.

22 Fig. 13 is a diagrammatic illustration of a cache manager with an object table
23 constructed according to yet another implementation.
24
25

1 Fig. 14 is a diagrammatic illustration of the sequence of operations from Fig. 10,
2 which shows a write optimizing technique for logging operations and objects affected by
3 write operations.

4 Fig. 15 is a diagrammatic illustration showing a sequence of logical application
5 operations and corresponding write graphs.

6 Fig. 16 is a diagrammatic illustration of the sequence of logical application
7 operations from Fig. 15, which shows the corresponding log records for those operations.

8 Fig. 17 is a diagrammatic illustration showing how a blind write operation
9 initiated by the cache manager affects a multi-object write graph.

10 Fig. 18 is a diagrammatic illustration of a cache manager with an object table that
11 is constructed to track dependencies introduced through both read and write operations.

12 Fig. 19 is a diagrammatic illustration showing a read operation, its corresponding
13 representation in terms of a write graph, and how the cache manager tracks any
14 dependencies in the object table.

15 Fig. 20 is a diagrammatic illustration showing a write operation, its corresponding
16 representation in terms of a write graph, and how the cache manager tracks any
17 dependencies in the object table.

18 Fig. 21 is a diagrammatic illustration showing a write graph with a combined
19 node formed from two collapsed nodes, and how the cache manager tracks this event.

20 Fig. 22 is a diagrammatic illustration showing a blind write operation to break a
21 cycle dependency, its corresponding write graph, and how the blind write affects the
22 object table.

23 Fig. 23 is a diagrammatic illustration showing how flushing an application object
24 affects the write graph and object table.

25

Fig. 24 is a diagrammatic illustration showing an excerpt of a stable log having log records and a conventional approach to identifying a point in the log to begin replaying operations during recovery.

Fig. 25 is a diagrammatic illustration showing the stable log of Fig. 24 and a recovery optimization technique for identifying a point in the log to begin replaying operations during recovery according to an aspect of this invention.

Fig. 26 is a diagrammatic illustration of a cache manager with an object table that is modified to track the starting log record for use in recovery.

Fig. 27 is a diagrammatic illustration showing the stable log having log records for a short-lived application object. Fig. 27 illustrates advancing the point to begin recovery according to the recovery optimization techniques.

Figs. 28A and 28B are exemplary write graphs produced by a sequence of operations that do not use a directed write-write edge of the present invention.

Fig. 28C is an exemplary write graph produced by a sequence of operations in accordance with the present invention.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

This invention concerns a recovery scheme that renders both data records and application programs persistent across system crashes. In general, the recovery scheme extends page-oriented, database style recovery to application programs. An application program's state is manifested in the application's address space. According to an aspect of this invention, the application state is treated as a single cached object, akin to a single memory page, which can be atomically flushed to a stable database. Application executions occurring between resource manager interactions are mapped to loggable operations that are posted to a stable log. The results of the application executions as

they impact other objects, such as data pages, are also captured as logged operations. The results of these operations are also from time to time flushed to the stable database. As a result, the recovery scheme allows integration of application recovery with database recovery.

The application recovery scheme is based on application replay. Application executions are logged during normal operation and are replayed during recovery. This reduces the recovery overhead for normal system operation while shifting more of the burden to the recovery process, wherein the logged application operations will need to be re-executed during recovery.

Fig. 3 shows a database computer system 50 having a computing unit 52, with a processing unit 54 and a volatile main memory 56, and a non-volatile memory 58 interfaced with the computer unit 52. The volatile main memory 56 is not persistent across system crashes. It is presumed to lose all data that it presently stores when a crash occurs. Main memory 56 can be implemented, for example, as volatile RAM. On the other hand, the persistent memory 58 is presumed to persist across a system crash. Examples of persistent memory 58 include disk arrays, disk drives (e.g., hard and floppy), read/write CD ROMs, tape backups, reel-to-reel, and the like.

The database computer system 50 is shown in an operational state in which one or more applications 60 are loaded in main memory 56 for execution on the processing unit 54. The application programs 60 are permanently stored on non-volatile memory (such as the persistent memory 58) and loaded into the main memory 56 when launched. The applications are representative of single threaded or multi-threaded applications. For purposes of continuing discussion, suppose that one of the applications is a long running application such as those that characterize workflow systems.

1 The main memory 56 further includes a resource manager 62 which maintains
2 temporary copies of the data pages and application states. The resource manager is
3 responsible for managing when to flush data objects and application objects, and hence
4 when to install operations into the persistent memory 58. It is also responsible for
5 posting operations from the volatile log to the stable log. This must be done before the
6 results of an operation are installed in the stable state, thus enforcing a write-ahead log
7 protocol. The resource manager 62 is callable by the application programs 60 and
8 mediates all data communication directed to and originating from the applications, as is
9 described below in more detail with respect to Fig. 5.

10 The resource manager 62 includes a volatile cache 64, a cache manager 66, a
11 volatile log 68, a log manager 70, and a recovery manager 71. The volatile cache 64
12 contains cached states of any executing application 60, and the data pages retrieved from
13 the persistent memory 58. The volatile log 68 tracks the operations performed by the
14 computer system.

15 The non-volatile memory 58 includes a stable database 72 and a stable log 74.
16 The stable database 72 maintains stable versions of the application address space and data
17 objects, and the stable log 74 maintains a stable sequence of logged computer operations.
18 The database 72 and log 74 are shown separately, but can be implemented in the same
19 storage subsystem.

20 The cache manager 66 manages the volatile cache 64 and is responsible for
21 retrieving data records from the stable database 62 and periodically flushing modified
22 data records back to the stable database 72. Additionally, the cache manager 66 manages
23 when to flush cached objects, including the application state as an object to be updated in
24 the stable database 72. The log manager 70 manages the volatile log 68 and facilitates
25

posting operations from volatile log 68 onto the stable log 74. In doing that, it enforces the write-ahead log protocol as directed by the cache manager 66.

The database computer system 50 is representative of many diverse implementations, including a database server for a network of PCs or workstations, an online server for Internet service providers, a mainframe computing system, and the like. The database computer system 50 runs an operating system (not shown), which is preferably a multitasking operating system which allows simultaneous execution of multiple applications or multiple threads of one or more applications. Examples of suitable operating systems include a Windows® brand operating system sold by Microsoft Corporation, such as the Windows NT® workstation operating system, as well as UNIX based operating systems.

One aspect of this invention is to make the applications 60 persist across system crashes, without requiring the applications to take steps to ensure their persistence. The recovery procedures implemented on the database computer system 50 are designed to work with conventional applications, which are not specially modified to account for, or even be aware of, recovery considerations. The applications are treated as individual objects that are flushed from time to time to the stable database 72. In this manner, application recovery can be likened to page-oriented database style recovery in that the monolithic application state is similar to a single database page.

To realize application recovery using page-like recovery technology, the system architecture of computer system 50 is designed to handle applications as individual, monolithic objects that can be independently flushed. The basic architecture involves two design issues: (1) how to atomically flush an operation consistent application state (which can be very large) as a single object, and (2) how to map application executions to

1 logical operations which change application state and can be posted to a stable log so that
2 the operations can be replayed during recovery.

3 Beyond this general architecture, however, are several optimizing features that
4 can be implemented to improve the efficiency and effectiveness of the application
5 recovery system. These other features include a modified cache manager that handles
6 such considerations as when to flush cached objects so as to avoid overwriting previous
7 states that might still be needed.

8 The following discussion first addresses the basic architecture, and then follows
9 with a description of the optimizing features.

10 11 Operation Consistent Application State

12 An object's operation consistent state is the state as it exists between operations.
13 The computer system 50 flushes operation consistent objects so that recovery, which
14 either re-executes an operation or bypasses it, works correctly. Database pages, when
15 flushed, are operation consistent. Page updates are short duration and under the control
16 of the resource manager; hence, operation consistency is achieved inexpensively with
17 standard techniques, e.g. latching or pinning.

18 Application state operation consistency is harder to provide. Applications execute
19 asynchronous to the resource manager. According to an aspect of this invention, the
20 application operations capture the application execution as state transitions between
21 interactions of the application with the resource manager. This aspect is described below
22 in more detail. A difficulty that arises is that the operation consistent application state as
23 of the last interaction with the resource manager no longer exists, and the cache manager
24 has no way of knowing when the application will again interact with the resource
25 manager to produce the next operation consistent application state.

1 There are several ways to provide operation consistent application state. One
2 technique is to capture and preserve the application state as of the most recent interaction.
3 Since the application state can be very large, capturing and preserving the entire state can
4 be expensive. However, this technique is a viable implementation and suitable for
5 recovery purposes, as large application states are capable of being atomically flushed to
6 stable storage using a conventional technique known as "shadowing," which is described
7 below.

8 Another technique is to force an application interaction with the resource
9 manager. The interrupted state of the executing application becomes operation consistent
10 by defining and logging the operations that precede and follow this state. To
11 demonstrate, suppose that the application state for application A is between interactions
12 with the resource manager during an application execute operation $Ex(A_i)$. The notation
13 " A_i " is used throughout this disclosure to refer to an application having an identifier "A"
14 taken at a state with a state ID of "i." To flush the application state at this intermediate
15 point, execution of the operation $Ex(A_i)$ is halted and the resulting intermediate state is
16 labeled A_{i+1} . The system defines and immediately flushes to the stable log a specially
17 marked execution operation $Ex'(A_{i+1})$, indicating a state transition from the interrupted
18 state A_{i+1} to the state as of the next interaction, i.e. A_{i+2} . The forced operation $Ex'(A_{i+1})$
19 makes the application state A_{i+1} operation consistent. Application state A_{i+1} can then be
20 flushed.

21 Three alternatives exist for replaying the operation $Ex'(A_{i+1})$ during recovery,
22 depending on when a crash occurs. When application A's persistent state identifier is:

- 23 1. Greater than $i+1$, operation $Ex'(A_{i+1})$ is bypassed like any other installed
24 operation;

2. Equal to $i+1$, operation $Ex'(A_{i+1})$ is replayed like any other application execute operation; or

3. Less than $i+1$, operations $Ex(A)$ are replayed normally through state i . Operation $Ex(A_i)$ is then replayed. Recovery bypasses the operation $Ex'(A_{i+1})$ following normal replay of $Ex(A_i)$ and simply increments application A 's state identifier to $i+2$. Replay of operation $Ex'(A_{i+1})$ can be avoided because replay of the preceding operation $Ex(A_i)$ at recovery (and hence, when $Ex(A_i)$ is not interrupted) inherently includes the execution of operation $Ex'(A_{i+1})$. This third case only arises when the system crashes between the log flush of forced operation $Ex'(A_{i+1})$ and the state flush of application state A_{i+1} .

Atomic Flush of Operation Consistent Application State

As part of application recovery, the database computer system 50 treats each executing application as a single object, which can be flushed from time to time to stable state in order to preserve snapshots of the application's address space. The database computer system 50 flushes the application state (which can be quite large) in an atomic operation.

Fig. 4 shows a portion of the database computer system 50 to illustrate a technique for atomically flushing the application state as a single object. The technique is known as "shadowing." The cache manager 66 maintains two versions of the application state: a current application state 80 kept in cache 64 and a lagging application state 82 kept in stable database 72. The lagging version 82 is the most recent version of the application state that has been flushed to the stable database 72. When the cache manager 66 decides to flush the current cached version 80 of the large application state, the cache manager 66 first writes the current cached version 80 to the stable database to form a

copy 80'. When the entire current version of application object has been written to the stable medium 72, the cache manager 66 moves a pointer (as represented by arrow 84) from the lagging version 82 to the new updated version 80' to place it logically within the stable database. Since the pointer 84 is small, it can be changed with a single page write. This enables the pointer to be moved between the two versions in an atomic operation. The earlier version 82 can then be discarded or overwritten.

Mapping Application Executions to Logical Loggable Operations

To ensure that operations are replayable during recovery, the operations are atomic and deterministic. An operation is said to be "atomic" if the external world that the operation sees during its execution appears to be constant, and the external world does not see the results of the execution until the operation completes. The operations are said to be "serializable" in that their execution is equivalent to an execution of the operations one at a time. An operation is said to be "deterministic" if, given the same system state as input, the result of execution against this state will always be the same output state.

To satisfy the atomic and deterministic criteria, all interactions between an application 60 and the external world (e.g., an end user, a database, a file, another application, etc.) are mediated by the resource manager 62. In this manner, the application is treated as a black box whose internal changes are not visible to the external world. These internal changes are not captured nor recorded in the volatile log. The application address space is intermittently exposed or impacted, however, every time the application interacts with the external world via the resource manager 62. Interactions with the resource manager thereby give rise to loggable operations that reflect different transitions between application states as the application executes. The application state transformations between interaction are hence logged as operations in the volatile log 68.

1 At recovery, these logged state transformation operations are replayed, with the affect
2 being that the hidden internal changes leading to each logged state are repeated.

3 Fig. 5 shows the system architecture of the database computer system 50 in more
4 detail. Individual application programs 60(1)-60(N) are executing on the computer. The
5 resource manager 62 provides an interface layer between each application and the
6 external world. In this illustration, the resource manager 62 mediates all communication
7 to and from the applications 60(1)-60(N) with respect to an end user at a terminal 86, a
8 data file in the cache 64, or another application. To interact with any external component,
9 an application calls to the resource manager 62 and the resource manager 62 facilitates
10 the task requested by the application. It does this by logging the application operation(s)
11 and then calling the requested system service that performs the requested task. This
12 intervening resource manager layer is said to “wrap” the requested task.

13 Execution of an application 60 is characterized as a series of loggable atomic
14 operations whose replay can recover the application. To capture application execution as
15 a series of loggable operations, the computer system 50 treats the code execution between
16 calls in the application as the log operation. Said another way, the resource manager 62
17 logs the operations as if it were calling the application, rather than the application calling
18 to the resource manager. This change in perspective results in an application operation
19 being “called” via a return from the resource manager 62 to the application 60. The
20 application operation “returns” to the resource manager via the application’s next call.

21 Given this shift in perspective, application execution is mapped into one of five
22 logical operations that are loggable in the volatile log 68. The five logical operations are
23 execute, initiate, terminate, read, and write.

24 1. Execute: A call from an application 60 to the resource manager 62 is treated
25 by the system 50 as a return from an application operation. A return to the application 60

1 from the resource manager 62 is treated as a call to an application operation. The
2 application execution between these interactions with the resource manager (i.e., starting
3 at a return from the resource manager and ending at the next call from the application to
4 the resource manager) is mapped to an *execute* operation.

5 Fig. 6 shows the logical execute operation. Suppose that the application is at a
6 state A_1 following a return from the resource manager. The application executes
7 instructions internal to the application, whose effects are hidden from the external world.
8 This execution transforms the application from a state A_1 to a state A_2 . Following this
9 execution, the application calls to the resource manager. The resource manager logs the
10 application Execute operation $Ex(A_1)$ denoting the transformation of application A from
11 state A_1 to state A_2 to the volatile log for subsequent posting by the log manager into the
12 stable log. As denoted in Fig. 6, the resource manager logs the application identifier A,
13 its state ID 2, and the execute operation Ex that resulted in the application state A_2 .

14 2. Initiate: This logical operation represents the application's first state transition
15 prior to the initial call to the resource manager 62. The resource manager 62 is notified
16 when the application is launched. The initial application state, e.g. its static code and data
17 structures, is read from stable memory during the launch. This action is mapped to a
18 loggable *initiate* operation. The initiate operation ends when the resource manager makes
19 the initial invocation of the application. The resource manager logs the $In(A)$ to the
20 volatile log for subsequent posting to the stable log.

21 3. Terminate: The *terminate* logical operation represents the application's final
22 call to the resource manager, instructing the resource manager to terminate the
23 application. This final application state transformation generates a "final state" for the
24 application that can be written back to the stable memory. When control returns to the
25 application, the application is expected to terminate cleanly and free up its resources. It is

not expected to call the resource manager again. The resource manager logs the Terminate(A) operation to the volatile log for subsequent posting to the stable log.

4. Read: The application 60 calls the resource manager 62 to read from an external system state, such as from a database page, perhaps in the cache 64. The resource manager 62 performs the read task, constructs a log record for this as a read operation that includes in the logged information the data values read and sufficient information so that the data read can be moved to the appropriate input buffers of the application state. The data is then moved to the application's input buffers and the log record is posted to the volatile log 68 and subsequently to the stable log. The return parameters of the read (i.e. the parameters that do not modify application state until control is returned to the application) become part of the log record for the next *execute* operation.

Fig. 7 shows a logical read operation following the execute operation described above with respect to Fig. 6. Suppose that the call made by the application to the resource manager at state A_2 is a call for a read task. The resource manager performs the read task and returns the values read from the object to the application. This return creates a change in application state to state A_3 . The resource manager logs the application identifier A and state identifier 3, the value of object O_1 , and the read operation R resulting in the application state A_3 . Thereafter, the log manager writes this log record to the volatile log and subsequently posts it into the stable log.

5. Write: The application 60 calls the resource manager 62 to write to external system state, such as to a database page that might already be in a buffer in cache 64. The resource manager 64 performs the write, logs the values written O.Val and the identity of the object O written in the log record in the volatile log 68. Any return parameters become part of the log record for the following *execute* operation.

Fig. 8 shows a write operation following the execute operation described above with respect to Fig. 6. Suppose that the call made by the application to the resource manager at state A_2 is a call for a write task. The resource manager performs the write task, logs the object identity O , its state ID 2, the values written O_2 , and the write operation W that results in the object state O_2 . The resource manager then returns any parameters resulting from the write task to the application. These return parameters are part of the input to the next execute operation.

One benefit of mapping the application execution into loggable operations is that these operations can be expressed entirely in terms of the application states. For the execute operation, for example, the application begins in one state and is transformed to another state by internal executions of the application. To the outside world, the execute operation can therefore be expressed as reading a first application state before the internal executions, and writing a transformed application state resulting from the internal executions. Table 1 shows the application operations characterized in terms of application states.

Table 1

Logical Operation	Expressed as Read/Write of Application State
Execute $Ex(A)$	Read application state, write application state.
Initiate $In(A)$	Write application state from the static state retrieved from stable memory. This writes the application invocation state instance.
Terminate $T(A)$	Write final application state.
Read $R(A)$	Read application state, write application state with read object data values that are included in the read log record.
Write	Writes do not effect application state. However,

W(O)	an application write transforms the written object from one state to another by overwriting its prior value with the after-image value stored in the write log record. Accordingly, a write operation writes data object state.
------	---

It is noted that there may be interactions that cannot be mapped into these five operations. For example, reading a message may consume the message as well; i.e. the application writes to the message queue by removing the message. This interaction is both a read and a write that cannot be optimized as above.

Fig. 9 shows an example series of loggable operations that are mapped from application executions. The loggable operations are designated by a circle: the legend “Int” within a circle stands for an initiate operation; the legend “Ex” within a circle represents an execute operation; the legend “R” within a circle stands for a read operation; the legend “W” within a circle represents a write operation; and the legend “T” within a circle stands for a terminate operation.

The initiate operation 90 writes the initial application state A_1 . The resource manager includes in a single log record an application identity A, its state ID 1, and the name of the operation Int. The log record is written in the volatile log and subsequently posted to the stable log.

An execute operation 92 reads the application state A_1 , performs some internal executions, and writes the application state A_2 by means of the application executing beginning in at state A_1 and the execution resulting in state A_2 . The resource manager logs the application identifier A, a state ID 2, and the execution operation Ex that resulted in the application state A_2 .

A read operation 94 reads the application state A_2 and an object O_1 . As above, the shorthand notation “ O_1 ” means an object with an identifier O taken at a state ID “1.” The

1 object value O_1 is read into the application buffers and results in a next application state
2 A_3 . The resource manager logs the application identifier A , its state ID 3, and the read
3 operation R that resulted in the application state A_3 . In addition, the resource manager
4 includes the object value O_1 in the log record. Writing the values read from the object
5 into the log record ensures that the values are available for redo of the application
6 operations during recovery in the event that the object O has been subsequently updated
7 and a subsequent value flushed to the stable database.

8 Unfortunately, in some cases, the values read from the object O can be large and
9 hence logging the entire object value is not desirable. Moreover, the log record
10 containing the object values is separate from, and often duplicative of, the data pages
11 holding the object O_1 which are occasionally flushed to the stable database. The system
12 and methods described herein address this problem by optimizing the read operation to
13 reduce the amount of data placed on the log. This optimization involves development of
14 a new cache manager, a topic that is discussed below with reference to Figs. 10-14 in
15 more detail.

16 An execute operation 96 transforms the application state from state A_3 to state A_4 .
17 The resource manager logs the application identifier A , a state ID 4, and the execution
18 operation Ex that resulted in the application state A_4 .

19 A write operation 98 writes a modified version of the previously read object,
20 designated as O_2 . The resource manager logs the object identifier O , its state ID 2, the
21 value O_2 written, and the write operation W that resulted in object state O_2 . This ensures
22 that the write parameters are available on the log for redo of the application operations
23 during recovery in the event that the object O_2 is not flushed to the stable database.

24 Similar to the read case, the value O_2 can be large and duplicated elsewhere in the
25 system, and thus logging the entire object value is not desirable. The system and methods

described herein address this problem by optimizing the write operation to avoid logging the value of O, by logging the application state that provided the data value for O. This write optimization involves development of a new cache manager, a topic that is discussed below with reference to Figs. 15-23 in more detail.

An execute operation 100 transforms the application state from state A_4 to state A_5 . The resource manager logs the application identifier A, a state ID 5, and the execution operation Ex that resulted in the application state A_5 .

A terminate operation 102 writes the final application state A_6 . The resource manager writes in a log record the application identifier A, a state ID 6, and the termination operation T that resulted in the application state A_6 .

The changes to the application during these operations are accumulated in the application state stored in the volatile cache. From time to time, the cache manager flushes the application state to stable storage. The flushed application state is tagged with a state ID. The flushing of the application state effectively installs all application operations which have been logged in the stable log that have a state ID less than the state ID of the flushed application state.

General Recovery

Following a system failure, the database computer system invokes a recovery manager 71 to recover the data pages (and other data objects) and application state lost during the crash. During redo recovery, the recovery manager 71 retrieves the most recently flushed data objects and application objects in the stable database and replays the operations in the log against the stable objects. The recovery manager 71 can be implemented as a conventional recovery manager which replays the stable log, beginning at a point known to be earlier than the oldest logged operation that was not yet installed.

1 The recovery manager compares the state ID of each logged operation in the stable log
2 with the state ID of a retrieved data object or application object. If the state ID of the
3 logged operation is later than the state ID of the stable object, the recovery manager
4 redoes that logged operation. This redo process returns the database computer system to
5 the previous state in which it was operating immediately prior to the crash, including the
6 recovered applications.

7 Another aspect of this invention involves techniques to optimize recovery to avoid
8 replaying operations that are rendered obsolete by subsequent operations. In this case,
9 the recovery manager is implemented to handle the recovery optimization techniques, as
10 is described in more detail below with reference to Fig. 24-27.

11 12 Read Optimization

13 In the recovery scheme described above, the read operation involves writing all of
14 the contents read from the object to the stable log in association with the read operation.
15 The logged operation can then described as reading and writing application state. This
16 type of operation, in which only a single object is written, and at most that object is read,
17 is referred to as a “physiological operation.” These operations are useful in that using
18 only such operations, recovery can be implemented using conventional cache managers
19 and cache management techniques. The cache manager need not be concerned about
20 object flushing sequence or preserving a certain object state because any data value
21 obtained from an object which was read, and hence which is needed to redo an
22 application operation is available directly from the stable log.

23 The benefits accruing to cache management as a result of logging only
24 physiological operations come at a cost. Treating an application read as a physiological
25 operation requires writing data, and often large amounts of data, to the stable log. This

1 reduces efficiency in the logging process and consumes I/O resources. Moreover, the
2 data written to the stable log is a copy of data in an object, which is maintained in volatile
3 cache and occasionally flushed to the stable database. It is wasteful to duplicate large
4 data objects in log records when these objects are available elsewhere.

5 Accordingly, an aspect of this invention is to optimize the logged read operation
6 to avoid writing the object's data to the log record. Generally, the optimizing technique
7 eliminates logging the read values by substituting, for the read values, names of the
8 objects from where the values are read in the log record. That is, rather than logging the
9 object value that is read, the read optimization technique involves logging the identity of
10 the object that is the source of the values being read. We call this a "logical read" and
11 denote it by $R(A,O)$, indicating that application A reads data object O for the input value
12 needed to transform application A's state; it does not get this input value from the log
13 record. For instance, a log record for the logical read operation includes the application
14 object's identifier A, its state ID, A.SID, the data object's identifier O, the data object's
15 state ID, O.SID, and an indication that a read operation was performed:

16
17 $\langle \underline{A}, A.SID, O, O.SID, R \rangle$
18

19 Other information may also be included, such as an index to a specific value set
20 contained in the object. Posting information that names the source of a data value, rather
21 than the value itself, substantially reduces the amount of information placed on the stable
22 log. When redoing a logged operation during recovery, the recovery manager 71 uses the
23 object name to locate the object and reads the value from that object.

24 Unfortunately, substituting object names for the actual values comes at a cost of
25 introducing dependencies between the objects in the cache. Attention must now be paid

1 to the order in which objects are flushed to stable storage. If objects are flushed out of
2 proper sequence, a particular state of an object may be irretrievably lost. An object name
3 contained in a logged operation would not enable restoration of the object values needed
4 by the operation if the data value for the object is not the same as the value that was
5 originally read from the object during normal execution.

6 Fig. 10 illustrates the dependency issue introduced by the read optimization
7 technique. Fig. 10 shows a sequence of operations comprising a read operation 110, an
8 execute operation 112, a write operation 114, and an execute operation 116. These
9 operations are identical to the operations 94-100 in Fig. 9. However, unlike the
10 procedure in Fig. 9, the value of the object that is read at operation 110 is not posted to
11 the log. Instead, only the object identifier O and state ID are posted. The object identifier
12 and state ID identify the exact data value needed by the logged operation.

13 The operation sequence in Fig. 10 introduces a dependency between the
14 application object A and the data object O. Assume, for example, that the cache manager
15 flushes the data object O to stable memory at state O_2 after the write operation 114
16 without having previously flushed the application object A to install the operations 110
17 and 112 preceding the write operation 114. Then, before the cache manager has an
18 opportunity to flush the application object A, the system crashes. Upon replay of the log,
19 the computer database system is unable to redo the operations to resurrect the true
20 application states A_2 - A_4 because the object state O_1 is not available. That is, the stable
21 database only contains the flushed object O at state 2, not at its initial state 1.

22 (Note that we do not describe the write 114 as reading application state A_3 .
23 Rather, write 114 is a physical write that gets the value written as O_2 from the log record.
24 This avoids additional flush dependencies.)
25

1 This dependency is explained in the context of an installation graph as a “read-
2 write edge.” That is, the write operation writes data into a read variable set which is read
3 in an operation preceding the write operation, thereby overwriting needed data to carry
4 out the read operation during recovery. Installation graphs and the read-write edge case
5 are described in detail in a publication by David B. Lomet and Mark R. Tuttle, entitled
6 “Redo Recovery after System Crashes,” Proceedings of the 21st VLDB Conference,
7 Zurich Switzerland, 1995. This publication is incorporated by reference.

8 To manage dependencies, the database computer system is equipped with a cache
9 manager that is attentive to flushing sequence. The cache manager is designed to ensure
10 that an application object is flushed to stable memory, thereby installing its operations,
11 before any modified data objects from which the application has read are flushed to stable
12 memory. The cache manager implements an object table which tracks active objects in
13 the volatile cache, and monitors flush order dependencies between those objects.

14 Fig. 11 shows a cache manager 120 with an object table 122. The object table 122
15 holds a list of objects that are presently stored, in the volatile cache or that have flush
16 dependencies with objects presently stored. These objects may be in the form of
17 application objects or data objects. Typically, the data objects have volatile (i.e. cache)
18 locations that are identified as memory pages. With regard to data objects, the object
19 table 122 is similar to prior art “page tables.” However, unlike prior art page tables, the
20 object table 122 also maintains a list of application objects, with each application object
21 comprising the application address space, and information with each entry that is used to
22 manage flush dependencies.

23 The object table 122 shows an entry 124 for the application object A and an entry
24 126 for the data object O which reflect respective object states following the read
25 operation 110. These entries contain information pertaining to the objects which is

1 organized in data structures 128 and 130. Each data structure has an object identifier
2 field 131, 132 to hold the object identifier (e.g., A or O), a state identifier field 133, 134
3 to hold the state ID for the value of the object, a dirty flag field 135, 136 which holds a
4 flag bit indicating whether or not the object has been modified in volatile cache without
5 those modifications being flushed to stable memory, and a cache location field 137, 138
6 to hold an address to a location in volatile cache where the current cached value of the
7 object physically resides. The data structure may further have a stable location field to
8 hold an address of the object in stable memory, although this field is not shown in this
9 example. Alternatively, the stable location may be derivable from the object identifier,
10 objectID, in field 131, 132.

11 Each data structure 128, 130 also has a predecessor field 139, 140 to hold
12 information for any predecessor object. An object is a "predecessor object" to a subject
13 object if that object must be flushed prior to flushing the subject object. The predecessor
14 field 139, 140 enables the object table 120 to track dependencies between the operations.
15 For the read operation, the dependency cases can be resolved into two rules: (1) only an
16 application object can be a predecessor; and (2) an application object has no predecessor.
17 The underlying reason for these rules can be better understood with a brief introduction to
18 a "write graph," which is a graph derived from an "installation graph," and is described in
19 the above incorporated article by Lomet and Tuttle.

20 Fig. 12 shows a write graph 144 for a read-write edge in which a read operation
21 reads a data object O at a first state during execution of the application object A, and
22 subsequently a write operation writes the data object O to create a second state of that
23 data object. In write graph notation, the circles represent nodes. A write graph node n is
24 characterized by a set of operations ops(n) and a set vars(n) of variables (i.e., objects)
25 written by the operations in ops(n). There is an edge between write graph nodes m and n

1 if there is an installation graph edge between an operation in ops(m) and an operation in
2 ops(n). The cache manager installs the operations of ops(n) by flushing the objects of
3 vars(n) atomically.

4 Write graph 144 has two nodes, an application node 146 with vars(146) = {A}
5 and a node 148 with vars(148) = {O}. The application node 146 shows that the read
6 operation has been performed which changes the application state (by reading values into
7 the application buffers) and that the application has continued its execution with an Ex(A)
8 operation. The data node 148 shows that the write operation affects the object state.

9 Write graph 144 demonstrates a flush order dependency between the application
10 object and data object. To ensure correct recovery of the application, the cache manager
11 flushes the application object represented by node 146, thereby installing the read
12 operation, prior to flushing the data object represented by node 148.

13 This write graph further illustrates that, for a logical read operation, an application
14 object A has no predecessor for which it is concerned. All paths between nodes 146 and
15 148 are at most a length of one. Only the data object O has a predecessor and that
16 predecessor is the application object A (which read it). The logical read operation, by
17 itself, thus reduces to a straightforward result. With reference again to Fig. 11, the
18 predecessor field 140 denotes a list of predecessors for the object O entry 130. The
19 predecessor entry shown contains the identifier for the application object A data record
20 128, denoted as the predecessor object PO. This predecessor is established when the read
21 operation 110 (Fig. 10) is encountered. The predecessor entry also includes a state
22 identifier for the originating object O, i.e., O.SID. That is, in the general case, an entry
23 on the predecessor list is represented as:

24
25 < O.SID, PO >

1
2 It is noted that a data object may have more than one predecessor. Hence, the
3 predecessor field 140 may contain a set of entries for multiple predecessor objects.

4 Since Fig. 11 illustrates a read operation, the application object has no
5 predecessor. As a result, the predecessor field 139 for the application A data structure
6 128 contains a null pointer, denoting the empty list.

7 Each data structure 128, 130 further includes a successor field 141, 142 to hold
8 information for any successor object. An object is a "successor object" of a subject
9 object if the subject object must be flushed before the successor object is flushed. The
10 successor field 141, 142 is primarily used as a bookkeeping function, to track successor
11 objects, as it adds no additional information that is not already contained in the
12 predecessor field. When flushing an object, the cache manager ensures that all real
13 predecessors are flushed beforehand. After flushing, the cache manager uses successors
14 only to clean up by removing the flushed object as a predecessor in other predecessor
15 lists. Less information is needed for successors, for example, object state ID, O.SID is
16 not needed. The cleanup is unconditional, taking place regardless of whether the
17 predecessor/successor is real or potential. It is noted, however, in an alternative
18 implementation, the successor field may be primarily relied upon, with the predecessor
19 field serving a secondary bookkeeping role.

20 The first statement of the read operation is that only an application object can be a
21 predecessor. The converse to this statement is that only an application object can have a
22 successor. In Fig. 11, the successor field 141 of the application A data structure 128
23 contains an entry for the object O data record 130. The successor entry is established
24 when the read operation 110 (Fig. 10) is encountered. The data object O has no
25

1 successor. As a result, the successor field 142 for the object O data structure 130 contains
2 a null pointer indicating an empty list.

3 Through the predecessor and successor fields in the object table, the cache
4 manager 120 tracks dependencies between the objects. When the cache manager 120
5 decides to flush an object to stable memory, the cache manager first checks the object
6 table 122, and particularly, the predecessor field of the object entry to determine whether
7 or not the object to be flushed has any predecessors. If a predecessor is listed for that
8 object, the cache manager will flush the predecessor object, assuming it is “real,” prior to
9 flushing the subject object.

10 The cache manager 120 distinguishes between “real” and “potential”
11 predecessors. A “real” predecessor object is one that has read an object whose state has
12 been changed by subsequent operations since the time the object was read by the
13 predecessor. A real predecessor must be flushed prior to the subject object to ensure
14 retention of a correct state in the stable database. In contrast, a “potential” predecessor
15 object is one that has read an object whose state has not changed since the time the object
16 was recorded as a predecessor. A potential predecessor does not require priority flushing,
17 although the cache manager continues to track potential predecessors because they may
18 turn into real predecessors. These are tracked by retaining object table entries for objects
19 with predecessors, even if they themselves are flushed and their values removed from the
20 cache.

21 Fig. 10 demonstrates the difference between real and potential predecessors. At
22 the read operation 110, the cache manager updates the predecessor list for the data object
23 O in the object table to reflect that the application object A is a predecessor. At this
24 point, however, application object A is only a “potential” predecessor because object O’s
25 value is still the same. Hence, application object A does not require flushing prior to the

1 data object O as the same application state can be recovered from re-reading data object
2 O.

3 However, at the write operation 114, the predecessor becomes a “real”
4 predecessor. Here, the data object O is modified by the write operation 114, thus
5 changing the state of O that the application object A read previously in the read operation
6 110. Now, application object A needs to be flushed prior to the data object O, or else
7 application object A will not be restored to the same application state during recovery
8 because the state 1 of data object O is irretrievably lost.

9 The cache manager determines whether a predecessor is “real” or “potential” by
10 comparing the current state identifier of the object to be flushed against the state identifier
11 of the same object as recorded in the entry of the predecessor list. For example, suppose
12 the cache manager 120 decides to flush data object O following the execute operation 112
13 (Fig. 10). The cache manager compares the current state ID of the data object O, which is
14 still state 1 at that point, with the state ID recorded in the entry for the predecessor
15 application object A contained in the predecessor field 140. In this case, object O’s state
16 ID in the entry is also 1. The state IDs match and thus, the application object A is only a
17 potential predecessor at this point. The cache manager is free to flush the data object O at
18 this point without first flushing application A. The predecessor entry for application
19 object A is maintained, however, in the predecessor field 140 of the object O entry 128.

20 Now, suppose that the cache manager decides to flush the data object O following
21 the write operation 114 (Fig. 10). The cache manager compares the current state ID of
22 the data object O, which is now state 2 following the write operation, with the state ID
23 recorded in the entry for the predecessor application object A contained in the
24 predecessor field 140. As before, the object state ID in the entry is 1. The state IDs no
25 longer match. Thus, the application object A has now become a real predecessor. When

1 faced with a real predecessor, the cache manager must first flush the predecessor, in this
2 case the application object A, prior to flushing the data object O. Flushing the application
3 object A effectively installs all of the operations (which in the example, all update
4 application A) through the write operation 114 (which accounts for the new object state
5 O_2).

6 Once the application object A is flushed, the predecessor entry contained in the
7 data object O's predecessor list 140 is removed. The cache manager deletes the
8 predecessor entry from the predecessor list 140. Since application object A may also be a
9 predecessor for other objects, the cache manager uses the application object A's successor
10 list 141 to inform any successor data objects (including data object O) that application
11 object A has been flushed and is no longer a predecessor to them.

12 When an application terminates, the cache manager scans the successor field 141
13 of the application object A to remove from the predecessor field of successor objects any
14 entries to the terminated application.

15 Fig. 13 shows an alternative construction of the object table. In Fig. 13, the object
16 table 150 contains an entry 152 for a data object O at state 1. This entry includes a data
17 structure 154 having an object identifier field 156, a dirty flag field 158, a cache location
18 field 160, a predecessor field 162, and a successor field 164. In data structure 154, the
19 predecessor field 162 contains an index to a separate predecessor table 166.

20 For each predecessor object, the predecessor field 162 contains a unique index to
21 an entry in the predecessor table 166 containing information used to identify and locate
22 the predecessor object. In this example, the entry in the predecessor table 166 contains a
23 real bit and an object identifier of the predecessor (i.e., $\text{objectID}_{\text{Pred}}=A$). The real bit
24 which is set (i.e., to a binary value 1) if the predecessor object is a "real" predecessor and
25 is reset (i.e., to a binary value 0) if the predecessor object is a "potential" predecessor.

1 When the cache manager decides to flush the data object O, the cache manager no longer
2 compares state IDs to determine whether a predecessor is real or potential. Instead, the
3 cache manager examines the real bit. If the real bit is set, the cache manager knows it
4 must flush the associated predecessor object before flushing the subject object. The
5 "real" bit is initialized to zero when an object O is read by an application. At the time
6 that the object O is subsequently written, all current potential predecessors (which have
7 real bit set to zero) have this bit set to one.

8 The read optimization techniques described in this section are beneficial because
9 they eliminate having to post the values obtained during a read operation onto the log.
10 Instead, the log only contains information to identify the object that was read. While this
11 reduced the amount of data to be logged, the read optimization techniques introduced
12 flush dependencies between objects. The cache manager thus keeps an object table which
13 tracks dependencies to ensure a proper flushing order.

14 15 Write Optimization

16 In the general recovery scheme introduced at the beginning of this detailed
17 disclosure, a write operation involves posting, to the stable log in association with the
18 write operation, all of the values that are written to an object. The logged operation can
19 be described as simply writing the object state of a data object. This yields a
20 physiological operation that can be handled using conventional cache managers and cache
21 management techniques. The conventional cache manager need not be concerned with
22 object flushing sequence or preserving a certain object state because any data value
23 written to an object, and hence is needed during recovery, is available directly from the
24 stable log.
25

1 data object state O_2 . Unlike Fig. 10, however, the value written to object O (i.e., O_2) at
2 the write operation 114 is not posted to the log. Instead, the cache manager logs the
3 identity of the data object that is written (i.e., O), the data object's state ID 2, the identity
4 of the application object A which is the source of the values written, object A's state ID
5 3, and the write operation W which results in object state O_2 . Posting these objects'
6 identities consumes much less memory and fewer I/O resources than posting the entire
7 value of the object state O_2 to the stable log. The application object identifier A and its
8 state ID identify the exact data value needed by the logged write operation.

9 The write optimization technique comes at the expense of introducing more flush
10 order dependencies to ensure proper installation of operations. In the read optimization
11 case described in the preceding section, flush order dependencies are comparatively easy
12 to handle. The dependency chain is at most one link in length. The application state in a
13 read dependency has no predecessors, and hence nothing ever needs to be flushed before
14 the application state itself. When the cache manager decides to flush an object, it flushes
15 all predecessor objects (i.e., any predecessor application objects) and then the subject
16 object. The read dependencies are thus "acyclic," meaning that each object can be
17 flushed atomically independently of other objects in a prescribed order, without requiring
18 the simultaneous atomic flushing of multiple objects.

19 Unfortunately, flush dependencies arising from write operations, when combined
20 with dependencies arising from read operations, can result in "cyclic" flush dependencies.
21 This means that an object that is both read and written by an application must be flushed
22 both before (actually, not later than) and after (actually, not earlier than) the application
23 object. Cyclic flush dependencies require atomically flushing both the data object and the
24 application object simultaneously, which presents significant complications.

Fig. 15 illustrates a cyclic dependency introduced by the write optimization technique. Fig. 15 shows a sequence of operations and how the operations are represented as write graphs. The sequence of logical operations includes a read operation 190, an execute operation 192, a write operation 194, an execute operation 196, a read operation 198, and a write operation 200.

Corresponding write graphs 202-212 are provided below each operation. The write graphs consist of nodes. Each node n identifies a set of uninstalled operations (i.e., the abbreviations above the dotted line within the nodes), denoted $ops(n)$, in correlation with a set of data or application objects written by the operations (i.e., the abbreviations below the dotted line within the nodes), denoted $vars(n)$. The cache manager usually sees the operations in serialization order. Including the operations in the write graphs in that order is fine because serialization is stronger than installation order.

At the read operation 190, the corresponding write graph 202 consists of a node containing application object A. The read operation 190 reads application state A_1 and data object state O_1 and writes application state A_2 . This is reflected in the write graph 202 as involving two nodes: one node containing the application object A and one node containing the data object O. The read operation is registered in the node containing the application object A because the operation writes the application state. The notation R_{190} (i.e., read operation 190) in the node containing the application object A indicates that the read operation writes object A. No operation is placed in the node containing object O, because the read operation does not write the object state.

When a new operation occurs, the operation is added to the write graph as follows:

- 1 1. Merge into a single node m all nodes n for which $\text{vars}(n)$ intersect ($\text{write}(\text{Op})$
2 intersect $\text{read}(\text{Op})$) is not null, where $\text{write}(\text{Op})$ is the set of variables written
3 by operation Op , and $\text{read}(\text{Op})$ is the set of variables read by Op .
- 4 2. If the resulting graph has a cycle, collapse each strongly connected region of
5 the graph into a single node. Each such node n has $\text{ops}(n)$ that equals the
6 union of $\text{ops}(p)$ of nodes p contained in its strongly connected region and
7 $\text{vars}(n)$ that equals the union of $\text{vars}(p)$.
- 8 3. For each node $p \neq m$, set $\text{vars}(p) = (\text{vars}(p) - \text{nx}(\text{Op}))$. This removes from
9 $\text{vars}(p)$ objects that become not exposed, where $\text{nx}(\text{Op}) = \text{write}(\text{Op}) -$
10 $\text{Read}(\text{Op})$.
- 11 4. Include a write-write edge so that unexposed objects that were removed from
12 $\text{vars}(p)$ are ordered to flush to the stable database after exposed objects
13 remaining in $\text{vars}(p)$ are flushed to the stable database.
- 14 5. Include a 'reverse' or 'inverse' write-read edge to ensure that objects in node
15 p are not exposed when they are flushed to install their operations. In other
16 words, an edge is defined from each node q to a node p , where the operation
17 from q reads the final version of the object in the node p . Previously, each
18 node p had node q as a potential predecessor.

19
20 The read operation 190 introduces a potential read-write edge in write graph 202
21 from the node containing A to the node containing O . This potential edge (shown as a
22 dashed arrow) indicates that a subsequent write or update of data object O to change its
23 state will create a real edge, thereby establishing a flush order dependency between
24 objects A and O . The direction of the arrow represents the flushing sequence in the flush
25 order dependency. The arrow points from the node containing object A to the node

1 containing object O (i.e., $A \rightarrow O$) to represent that the application object A must be
2 flushed before the data object O.

3 The execute operation 192 reads the application state A_2 and writes the application
4 state A_3 . The node containing the object A in the write graph 204 is expanded to include
5 the execute operation (i.e., Ex_{192}) because the execute operation 196 writes application
6 state A_3 . The node containing object O remains void of any operations.

7 The write operation 194 reads application state A_3 and writes the object state O_2 .
8 The write operation is reflected in the write graph 206 by placing the notation W_{194} (i.e.,
9 write operation 194) in the node containing the data object O. Notice that the write
10 operation 194 does not write the application state, and thus the write operation is not
11 added to the node containing application A.

12 The write graph 206 also shows a real read-write edge caused by the read and
13 write operations 190 and 194. That is, the previous potential edge has now been
14 converted to a real edge by virtue of the sequence of read-write operations 190 and 194.
15 This read-write edge introduces a flush order dependency between application object A
16 and data object O. To ensure correct recovery of the application, the cache manager must
17 flush the application object A, thereby installing the read operation R_{190} , prior to flushing
18 the data object O. The read-write edge is indicated by a solid arrow, the direction of
19 which indicates the flushing sequence in the flush order dependency. Here, the
20 application object A must be flushed before the data object O and thus, the arrow points
21 from the node containing object A to the node containing object O (i.e., $A \rightarrow O$).

22 The write operation 194 also introduces a potential edge in write graph 206 from
23 the node containing O to the node containing A. This potential edge indicates that a
24 subsequent write or update of data object A to change its state will create a real edge,
25 thereby establishing a flush order dependency between objects A and O.

1 The execute operation 196 reads application state A_3 and writes application state
2 A_4 . Since the execute operation 196 writes application object A, the application node A
3 of the write graph 208 is expanded to include that operation (i.e., Ex_{196}). The execute
4 operation 196 does not write the data object state, and thus the execute operation is not
5 added to the node containing data object O.

6 The execute operation 196 introduces a real dependency between the data object
7 O and the application object A, as indicated by the write-execute edge. This dependency
8 arises because the data object state O_2 can only be regenerated from values found in the
9 output buffers at application state A_3 , which is about to change as a result of the execute
10 operation 196. Since the write optimization technique eliminates logging of the write
11 values to the stable log, the recovery manager must obtain those values from the output
12 buffers of application state A_3 to replay the write operation 194.

13 To ensure correct recovery of the data object O, the cache manager must flush the
14 data object O, thereby installing the write operation 194 which produces state O_2 , prior to
15 flushing the application object A. The write-execute edge is indicated by the solid arrow
16 pointing from the node containing O to the node containing A, thereby indicating an
17 $O \rightarrow A$ flushing sequence in the flush order dependency.

18 Unfortunately, the two dependencies between objects A and O are cyclic (i.e.,
19 $A \rightarrow O \rightarrow A$). As shown in the write graph 208, application object A must be installed
20 before data object O (i.e., $A \rightarrow O$) to ensure recovery of the application and the data object
21 O must be installed before the application object A (i.e., $O \rightarrow A$) to enable replay of the
22 write operation 194. This cycle can only be handled in full by flushing both objects A
23 and O simultaneously and atomically. This poses a problem.

24 To break such cycles, the cache manager 66 assumes an active role by timely
25 introducing "blind writes" that effectively preserve the state of data object on a log

1 record. In a blind write operation, the current value of the data object O is written to the
2 log in a manner similar to the general unoptimized write case discussed earlier in this
3 disclosure. The blind write leaves the value of data object O unchanged, but writes an
4 after-image of its value on the stable log. As a result, the data object O can be
5 regenerated from this log record, rather than relying on regeneration of a specific state of
6 the application object A.

7 Accordingly, the dependency cycle is broken. This enables an ordered flushing
8 sequence of first flushing the application object A and then flushing the data object O.
9 That is, once the cycle is broken, the cache manager can atomically flush objects one-by-
10 one, rather than having to flush multiple objects simultaneously and atomically.

11 The cache manager flushes the objects one-by-one according to a predetermined
12 acyclic flushing sequence. Preferably, as described above, the application object A is
13 flushed before the data object O is flushed. Thus, for example, a method in accordance
14 with the present invention comprises the step of flushing the data object O to the non-
15 volatile memory (i.e., the stable database) after the application object A has been flushed
16 to the non-volatile memory. However, according to another embodiment of the present
17 invention, the acyclic flushing sequence is arranged such that it is the application object
18 that is written to the log so that the data object O is flushed before the application object
19 A is flushed.

20 The way the cache manager identifies cycles and actively imposes blind writes is
21 best understood in the context of the write graphs. The process, as it pertains to write
22 graphs, involves three general steps. Also introduced is the "intermediate write graph,"
23 which is the graph formed before the cycles are collapsed.
24
25

- 1 1. Add each new operation to the intermediate write graph, either including it in
2 a node with existing operations or giving it a node of its own. The
3 intermediate write graph can have cycles.
- 4 2. Collapse nodes affected by cycles into a single node n (i.e. all intermediate
5 write graph nodes of the strongly connected region are collapsed into a single
6 write graph node.). The resulting node n has $\text{vars}(n)$ consisting of multiple
7 objects.
- 8 3. Remove all objects, but one, from the single node. This reduces $\text{vars}(n)$ to
9 containing a single object that needs to be flushed in order to install the
10 operations of the node n . The removal of objects can be accomplished
11 through normal write operations, or through a series of blind writes.

12
13 These three steps result in a new write graph containing nodes p with $\text{vars}(p)$
14 having a single variable that can be flushed by itself. The edges connecting these nodes
15 impose an order to the flushing of the objects, but the need to atomically flush multiple
16 objects is removed.

17 The edges that impose a flushing order are determined by a predetermined
18 acyclic flushing sequence. The acyclic flushing sequence is structured such that the
19 object that is removed from the cycle dependency by the blind write is flushed to the
20 stable database after the other object of the original cycle dependency. In other words,
21 the object that is not removed from the cycle dependency by the blind write is flushed to
22 the stable database before the object that is removed from the cycle dependency is flushed
23 to the stable database. The inventor has determined that such an acyclic flushing
24 sequence ensures recovery while providing effective cache management.

25

Step 1: Build the Intermediate Write Graph

The intermediate write graph is constructed by the cache manager 66 by performing the following steps for each operation:

1. Identify one or more objects that are both read and written by the operation., i.e. $\text{write(Op)} \cap \text{read(Op)}$.
2. Intersect the object(s) of step 1 with each set of existing objects associated with a present write graph node n , i.e. objects in $\text{vars}(n)$.
3. If all intersections are null, put the operation into its own node.
4. If an intersection is not null, merge all nodes with non-null intersections with the objects of step 1 into a single node.
5. Form edges between intermediate write graph nodes n and m based on when edges exist between the operations of $\text{ops}(m)$ and $\text{ops}(n)$ in the installation graph.
6. Remove the objects $\text{nx(Op)} = \text{write(Op)} - \text{read(Op)}$ from $\text{vars}(p)$ of any other node that currently contains them.

This method is repeated as new operations are executed and the intermediate write graph is built one operation at a time in operation execution order. A more detailed construction of one exemplary cache manager, and an object table which tracks write dependencies in a manner which effectively handles multi-object nodes and blind write strategies, is described below with reference to Figs. 18-23.

Step 2: Collapse Intermediate Write Graph Cycles

When a cycle is created, such as the cycle between the nodes containing A and O in the intermediate write graph 208 of Fig. 15, the affected nodes are collapsed into a single node. That is, all intermediate write graph nodes of a strongly connected region are collapsed into a single write graph node. Write graph 210 shows a combined node containing both objects A and O. This combined node contains the union of all operations and objects from the original two nodes. Collapsing intermediate write graph 208 results in the upper node of write graph 210. (The write graph is defined to be acyclic, while the intermediate write graph has cycles.)

Step 3: Reduce Objects In Node to One Object

Forming a combined node containing both A and O has not removed the dependency cycle; rather, both A and O must still be installed atomically together. To break the cycle so that variables can be flushed one by one, all but one object is removed from the node containing multiple objects. This can be done as a result of normal operation, or through a series of blind writes imposed by the cache manager.

With continuing reference to Fig. 15, the read operation 198 involves reading both the data object state O_2 and a new application state B_1 , and writing application state B_2 . The read operation 198 is reflected in the write graph 210 by addition of a node to contain object B and the inclusion of R_{198} (i.e., read operation 198) in that node. Additionally, the read operation 198 introduces a potential read-write edge from the node containing B to the node containing A, O. This potential edge indicates that a subsequent write or update of data object O to change its state will establish a flush order dependency between objects B and O in which the read operation 198 must be installed (by flushing object B) before installation of the operations 190-196.

1 The write operation 200 reads the application state A_4 and writes object state O_3 .
2 The corresponding write graph 212 is expanded to include a third node which contains
3 object O and W_{200} (i.e., write operation 200). This operation does not join the existing
4 node containing A, O because $\text{write}(200) \cap \text{read}(200)$ is null. The potential read-
5 write edge becomes a real "inverse write-read" edge as a result of this write operation
6 200. The read operation 198 (R_{198}) has read the last version of O written by write
7 operation 194 (W_{194}). This means that a real flush order dependency now exists because
8 data object O 's state has been changed in the write operation 200. The flush order
9 dependency dictates that the operation 198 in the node containing object B must be
10 installed prior to the operations 190-196 in the node containing objects A, O . A second
11 flush order dependency is also created by a read-write edge resulting from the write
12 operation. In this dependency, the application object B must be flushed, thereby
13 installing the read operation 198, prior to flushing the data object O .

14 The purpose of the inverse write-read edge is to ensure that data object O is not
15 exposed when the node with operations 190-196 has no predecessors. This permits the
16 operations 190-196 to be installed by flushing only A .

17 Notice that the result of write operation 200 removes data object O from the node
18 containing operations 190-196. An object can only reside in one write graph node, which
19 is the last node to write the object. Data object O is in $\text{nx}(200)$ and hence is removed
20 from the node containing operations 190-196. Here, the node containing write operation
21 200 is the last node to write object O , and hence, data object O resides only in that node.
22 No subsequent operation can remove it from that node without also writing it. Because
23 W_{194} and W_{200} both write data object O , and replay of W_{194} does not guarantee the ability
24 to replay W_{200} , there is an installation edge from W_{194} to W_{200} . This edge results in a write
25 graph edge from the node with operations 190-196 to the node with operation 200. There

1 is also an edge from R_{190} to W_{200} so this is a case where a write graph edge results from
2 two installation graph edges.

3 This is a case in which an object is removed from a multi-object node as a result
4 of normal operation. As a result of the write operation 200, the dependency cycle that
5 existed in the intermediate write graph 208 is now broken. That is, a single object A can
6 now be flushed to install all operations 190-196 in the node, including the write operation
7 194 that originally affected the data object O.

8 In terms of the write graph, the write operation renders the data object O
9 “unexposed” in the collapsed node of the write graph 212. An “unexposed” object of a
10 write graph node is one that has a write operation for it in a succeeding node and no read
11 operations following the current node that also do not follow the succeeding write. As a
12 result, an unexposed object does not need to be flushed in order to install the operations
13 in the preceding node that wrote that object as no succeeding operation needs the value
14 that it wrote. Conversely, an “exposed” object in a node is an object that needs to be
15 flushed to install the operations in the node that wrote that object. In the Fig. 15 example,
16 the application object A is “exposed” in the collapsed node. Although an unexposed
17 object does not need to be flushed, it is still preferably flushed to the stable database after
18 the exposed object is flushed. In this manner, recovery is ensured while providing
19 effective cache management.

20 Fig. 16 shows the corresponding log records for the sequence of operations 190-
21 200 from Fig. 15. As a result of the log optimization technique, the log record for the
22 write operation 194 does not contain the value written to the data object O (i.e., O_2).
23 Instead, the log record for write operation 194 contains only the data object identifier O,
24 the data object O’s state ID 2, the application object identifier A, its state ID 3, and the
25 write operation W that resulted in data object state O_2 .

Fig. 16 also shows another technique for reducing the number of objects in a multi-object combined node. The cache manager may not wish to wait for a subsequent write operation of one of the objects in the write graph node, such as write operation 200, because such operations cannot be foreseen and are not guaranteed to occur. Accordingly, the cache manager can impose its own write of an object in the multi-object node. The cache manager performs a "blind identity" write which writes the value of the object onto the stable log. Fig. 16 shows a blind write operation 216 which writes the values of the data object O at state 3, i.e., O_3 , to the log record. The blind write creates an after-image of the data object O on the log. That is, the blind write in this case is an identity write because the identical value of data object O, which is the same at both states 2 and 3, is written to the log. The state ID is stepped from 2 to 3 to maintain the convention introduced earlier in this disclosure.

Once the value O_3 is posted to stable log and all nodes that precede the node with operations 190-196 have been installed, i.e. the node with R_{198} , the cache manager is free to flush the application object A, thereby installing operations 190-196. If the system crashes after object A is flushed and application state A_3 is irretrievably lost, subsequent operations involving the data object O at state 3, can be replayed using the values O_3 on the stable log, rather than the values from the output buffers of a regenerated application state A_3 . Blind writes come at a cost of writing larger amounts of data to the log, but this cost is minimal in comparison to the advantages gained by the write optimization techniques in which a high percentage of writes do not result in posting entire object values to the log.

Although data object O does not need to be flushed to the stable database because it is written to the stable log, it is still preferably flushed to the stable database, and more preferably, it is flushed to the stable database after the exposed object A is flushed. This

1 subsequent flush is used to manage the cache. That is, the object is flushed when it is
2 desired to drop the object value from the cache. This allows the object to be stored
3 somewhere other than the cache where it can be retrieved if it is needed to be read or
4 updated again in the future.

5 Therefore, in a database computer system having a non-volatile memory, a
6 volatile main memory, and an application object which executes from the main memory,
7 wherein the non-volatile memory includes a stable log, a computer-implemented method
8 in accordance with the present invention comprises the following steps: executing the
9 application object to perform operations which read data from, and write data to, a data
10 object; posting to the stable log a log record for each operation involving the reading or
11 writing of data, the log record containing a reference to either the application object or the
12 data object to identify that referenced object as a source for the data that is read from or
13 written to; establishing flush order dependencies between the application object and the
14 data object, wherein some of the flush order dependencies become cyclic indicating a
15 condition in which the application object should be flushed not later than the data object
16 and the data object should be flushed not later than the application object; detecting a
17 dependency cycle; and following detection of the dependency cycle, writing one of the
18 application object or the data object to the stable log to break the dependency cycle to
19 enable the application and data objects to be flushed sequentially according to an ordered
20 flushing sequence.

1 Preferably, the writing step comprises writing the data object to the stable log.
2 More preferably, the method comprises the step of flushing the application object to the
3 non-volatile memory after the data object is written to the stable log. More preferably,
4 the method comprises the step of flushing the data object to the non-volatile memory
5 after the application object has been flushed to the non-volatile memory. Alternatively,
6 one can write the application object to the stable log, and then flush the data object first to
7 non-volatile memory and then flush the application object to the non-volatile memory.

8 The cache manager-imposed blind write has the same affect of removing an object
9 from the collapsed node in the write graph as a normal write operation. But such a write
10 is under the control of the cache manager, and hence the cache manager can use such
11 writes to help it manage the cache.

12 Fig. 17 illustrates the effect of a blind write on the combined node in the write
13 graph 210 of Fig. 15. In a blind write, the cache manager posts the current value of the
14 data object O to the stable log. This is represented in a write graph 211 as a new node
15 containing the object O and a blind write operation (i.e., W_{216}). Since the value of O is
16 written to the log, the data object O does not need to be flushed concurrently with the
17 flushing of application object A, and hence the $O \rightarrow A$ dependency is removed. The blind
18 write thereby breaks the dependency cycle.

19 In write graph terms, the data object O is no longer "exposed" in the combined
20 node and is withdrawn from that node. The cache manager no longer needs to flush
21 object O as part of the installation of the operations 190-196 in the combined node
22 because it does not matter what object O's value is. The cache manager need only flush
23 the exposed application object A to install all operations in the node, including those that
24 had written data object O, even though data object O is not flushed. Preferably, however,
25

1 data object O is flushed to the stable database after the application object A in order to
2 provide effective cache management.

3 It is noted that, for combined nodes having more than two objects that require
4 simultaneous flushing, the cache manager blind writes all but one object to the stable log.

5 Although the description herein is directed to an application object A and a data
6 object O, and the effects of a blind write with respect thereto, it should be noted that a
7 blind write can be used to break up a node containing multiple data objects and/or
8 multiple application objects. An acyclic flushing sequence is used by the cache manager
9 after a blind write to a multi-object node such that the object that is removed from the
10 node, regardless of whether it is a data object or an application object, is flushed to the
11 stable database after the object that remains in the node, regardless of whether that object
12 is a data object or an application object, and regardless of whether the multi-object node
13 is the result of a cyclic flush dependency or whether it arose in some other manner.

14 In other words, in accordance with the present invention, in a database computer
15 system having a cache manager which occasionally flushes objects from a volatile main
16 memory to a non-volatile memory to preserve those objects in the event of a system
17 crash, and wherein a dependency cycle exists between at least two objects such that the
18 two objects should be flushed simultaneously, a computer-implemented method
19 comprises the following steps: detecting a dependency cycle; and writing one of the two
20 objects to the stable log to break the dependency cycle to enable the two objects to be
21 flushed to the non-volatile memory in a sequential manner according to an ordered
22 flushing sequence. The method preferably comprises the step of flushing the objects
23 according to the ordered flushing sequence after the writing step.

24 Figs. 28A and 28B are exemplary write graphs produced by a sequence of
25 operations that do not use a directed write-write edge of the present invention. Each of

1 these write graphs is produced using a different prior technology. Fig. 28C is an
2 exemplary write graph produced by a sequence of operations in accordance with the
3 present invention, as described above.

4 Given the sequence of operations: (1) $F(z) = \{x,y\}$ (the log operation reads z and
5 writes x and y ; (2) $G(x) = w$ (the log operation reads x and writes w); and (3) $H(\) = x$ (a
6 blind write of x with some data that is stored in the log record), one type of write graph is
7 shown in Fig. 28A. There is a read-write edge 405, similar to that described above with
8 respect to Fig. 12, between the two nodes 400 and 410. The read-write edge 405 exists
9 because H writes into x , which G reads. All the operations that write x are together in the
10 same node 410. Thus, there is an undesirable cycle dependency in the node 410, which
11 the systems and methods of the present invention break apart.

12 Fig. 28B shows another write graph. In Fig. 28B, the node 410 has been replaced
13 two nodes 420, 430 using the blind write operation H . An inverse write-read edge 415 is
14 formed between the nodes 400 and 420 because G reads x , and is preferably installed
15 before F to make x unexposed. The blind write of x has removed x from the node 420
16 containing F , even though F writes x . A read-write edge 425 is formed between the
17 nodes 400 and 430, because H writes into x , which G reads.

18 As shown in Fig. 28C, after the blind write, in accordance with the present
19 invention, a write-write edge 435 is provided from the node 420 to the node 430, because
20 both nodes write x . The write-write edge 435 provides effective cache management, as
21 described above, by ordering that the flushing of the object x takes place after the
22 flushing of the object y . Thus, the object that is written to the stable log by the blind
23 write (object x) is flushed after the other object in the original node (object y).

24 It should be noted that Figs. 28A – 28C do not illustrate a cyclic dependency that
25 has been collapsed. Instead, one operation writes two objects, requiring that the objects

1 be flushed atomically. In Fig. 28A, there is no way to subsequently flush the objects
2 separately. Once objects are part of an atomic flush set, they remain a part of it, and must
3 be flushed to the disk to ensure atomicity. In Fig. 28B, the write graph permits the
4 objects to be flushed separately when there is an appropriate write that makes one of the
5 objects not exposed. In Fig. 28C, the write graph does the separation of Fig. 28B, and
6 adds an additional flush order constraint, as indicated by the write-write edge 435
7 between the write graph nodes 420 and 430.

8 Fig. 18 shows a cache manager 220 with an object table 222, which are
9 configured to track dependencies, including cyclic dependencies, and to manage the
10 object flushing sequence to properly handle the dependencies. The object table 222 holds
11 a list of objects that are presently stored, and in some cases recently stored, in the volatile
12 cache. These objects may be application or data objects.

13 The object table 222 shows an entry 224 for the data object. The entry is
14 organized in data structure 226 having an object identifier field 228, a dirty flag field 230,
15 a cache location field 232, and a node field 234. The node field contains an index to a
16 separate node list 236 of intermediate write graph nodes. These nodes all write to the
17 object with entry 224. Given that operations write at most one object, an operation can
18 always be associated with exactly one entry in the object table, i.e. the entry whose object
19 it wrote. All intermediate write graph nodes also have operations that write exactly one
20 object. The node list is a list of these intermediate write graph nodes containing
21 operations that write the object table entry.

22 The node list 236 is a data structure containing various entries 1, ..., N. Each
23 entry contains a "Last" field 238 that holds data indicating the last update to the object O
24 as a result of the operations of the node. The "Last" field 238 is set to the state identifier
25 of the object at its last update by operations of the node described by node list entry 236.

1 The node list entry also has a node identifier 240 to identify the write graph node into
2 which this intermediate graph node has been collapsed should the node be part of a cycle
3 (a strongly connected region) in the intermediate write graph. In this implementation, the
4 node ID field 240 is an index to a separate node table 246. This data structure contains an
5 entry 248 for write graph nodes that are produced as a result of an intermediate graph
6 collapse. Each such write graph entry has a list of all intermediate graph nodes from
7 which it is constituted via a collapse. These intermediate write graph nodes are identified
8 by pairs $\langle O, O.sid \rangle$.

9 As explained above with reference to Fig. 15, an object can be written by
10 operations in more than one node. The write graph 212 (Fig. 15), for example, shows that
11 data object O, while only requiring flushing in one node, is updated by operations in two
12 different nodes. The node ID fields 240 of all intermediate write graph nodes are "null"
13 until a cycle exists. When a cycle arises, the node ID of the intermediate write graph
14 nodes in the cycle are set to the write graph node identified by entry 248 in the node table
15 246, which includes the intermediate nodes of the cycle.

16 Each node list entry in node list 236 further has a predecessor list 242 and a
17 successor list 244. These lists are similar to those described above with respect to Fig. 11
18 in that they reference predecessor or successor write nodes (in this case, intermediate
19 write graph nodes) which should be flushed before or after the subject node. Each item in
20 the predecessor list 242 or successor list 244 must identify a predecessor or successor
21 node. Since an object can be written by operations in multiple write graph nodes, the
22 object's identifier is no longer sufficient for this node identification. The node can be
23 identified, however, via a pair $\langle \text{object id}, \text{state id} \rangle$, where the state ID is that of the Last
24 attribute for the node at the time the write graph edge was formed. (This can be used in a
25

1 look up that finds the node with the smallest Last value that is greater than this state ID.)

2 Thus, a node on a predecessor or successor list can be represented by:

$$N_x = \langle \text{Object ID of } X, \text{State ID of } X \rangle$$

3
4
5
6 In addition, as in Fig. 11, real and potential predecessors need to be distinguished.
7 This is done by storing with the predecessor list entry the state ID of the current object O
8 of entry 226 that was read by the first operation causing the edge described by the
9 predecessor entry. The state ID is denoted by $\text{first}(N_x, O)$. Thus, a predecessor list entry
10 is represented by the following format:

$$\langle \text{first}(N_x, O), N_x \rangle$$

11
12
13
14 The node being referenced in the predecessor and successor lists is an
15 “intermediate node,” not the write graph node. Multiple intermediate nodes may
16 comprise a write graph node, which is found from the entries via the Node ID field 240.

17 A successor list entry need only identify the successor intermediate node by a pair
18 $\langle \text{object id, state id} \rangle$.

19 The entries 1-N in the node list 236 are ordered according to the update order
20 sequence. This sequence is readily derived in the data structure by placing the entries in
21 ascending order according to the state identifier in their “Last” field 238.

22 The cache manager 220 uses the object table 222 to track the flush order
23 dependencies that arise in both read and write operations. Consider the case of a read
24 operation. Fig. 19 shows the read operation 190 from Fig. 15 in more detail. Read
25 operation 190 involves reading both application state A_1 and object state O_1 , and writing

1 application state A_2 . The intermediate write graph fragment 202 for this operation
2 includes a node 250 containing A and the read operation R_{190} , and a node 252 containing
3 O without any operations. The read operation 190 results in a potential edge from the
4 node containing A to the node containing O, indicating that a subsequent write or update
5 of data object O to change its state will create a real edge.

6 As a result of the read operation, the cache manager creates a node entry 254 for
7 the data object O's node list 236 which recognizes object A as a predecessor. Entry 256
8 is only a "potential" node list entry at this point since a write graph node technically only
9 exists when uninstalled operations write into variables. That is, the node containing data
10 object O becomes a write graph node in write graph 206 following the write operation
11 194. A node is shown in Figs. 15 and 19 to help describe how data object O is handled.

12 More particularly, node list entry 256 has a "Last" field 238 set to "1," the state
13 ID of data object O's last update, and a node ID field set to "null", indicating that this
14 node has not taken part in a "collapse". The predecessor list 242 is updated to reference
15 the predecessor application object A. This node reference includes the predecessor object
16 ID "A," and A's state ID of 2. In addition, to determine when this edge is real or
17 potential, the node reference includes "first($\langle A, 2 \rangle, O$)," indicating the state ID of data
18 object O when first read by application object A in this node, which is 1. The edge is real
19 only if data object O has a state ID that is greater than 1. Nothing is placed in the
20 successor list 244.

21 Similarly, the cache manager creates a node entry 256 for the application object
22 A's node list which recognizes data object O as a successor. Entry 256 contains in its
23 "Last" field 238' the state ID of "2" for the application object A's last update and in the
24 node ID field 240' it contains the value null, indicating that this intermediate write graph
25 node is not part of a cycle and hence has not taken part in a collapse. The successor list

1 244' of entry 256 is updated to reference the successor data object O. This successor
2 reference to identify the node for object O includes the successor object ID "O," and O's
3 state ID of 1. Nothing is placed in the predecessor list 242'.

4 Next, consider the case of the write operation. Fig. 20 shows the write graph 208
5 following the execute operation 196 and write operation 194 from Fig. 15 in more detail.
6 Write operation 194 involves reading application state A_3 and writing object state O_2 .
7 Execute operation 196 involves reading application state A_3 and writing application state
8 A_4 . The write graph 208 as a result of this operation includes the node 251 containing A
9 and the node 253 containing O. The write operation 194 represented in node 253
10 changed the data object state from state O_1 to state O_2 , thereby changing the previous
11 "potential" edge to a "real" edge (as represented by the solid arrow). This correlates to
12 object A becoming a real predecessor to object O. Additionally, recall in the write graph
13 206 of Fig. 15, a second potential edge had been created as a result of the write operation
14 194 because data object O, to be replayed, must obtain values from application object A
15 at state 3. This successor edge becomes real in write graph 208 of Fig. 20 because the
16 downstream execute operation 196 changes the application state from state 3 to state 4.
17 Thus, the write graph 208 in Fig. 20 shows two real edges between the nodes 251 and
18 253.

19 The old entry 254 representing a potential write graph node for data object O is
20 replaced by a real write graph node list entry 262. Entry 262 for data object O is created
21 in response to the writing of data object O at operation 194. The entry 262 has a "Last"
22 field 238" set to the object O's state ID following the write operation 194 (i.e., state ID =
23 2), and a node ID field 240" set to null. The predecessor list 242" in entry 262 includes
24 the same reference to predecessor object A as is contained in the predecessor list 242 in
25 entry 254. The successor list 244" in entry 262 is updated to reference the successor

1 object A. This reference includes the successor object ID "A" and A's state ID of 3.
2 Whether a successor is considered "potential" or "real" has little impact. When the
3 predecessor is flushed, the predecessor is removed from its successors' predecessor list
4 entries, regardless of whether it is real or potential.

5 With respect to the node list entry 256 for application object A, the "Last" field
6 238' has been updated to reflect a state 4 since this is the state at the execute operation
7 196. (Figure 20 shows the data structures after this operation, but before the collapse of
8 the cycle that is now present.) The cache manager also updates the predecessor list 242'
9 of the node list entry 256 for application object A to reference the "potential" predecessor
10 object O. This node reference includes the predecessor object ID "O," and O's state ID of
11 2. In addition, to determine when this edge is real or potential, the node reference
12 includes "firststr(<O,2>,A)," indicating the state ID of application object A when first read
13 to write the data object O at state 2, which is a state ID of 3. The edge is real only if
14 application object A has a state ID that is greater than 3. Figure 20, because it shows the
15 write graph 208 after the execute operation Ex_{196} , shows the edge as real, with the Last
16 field of 256 set to 4.

17 Notice that the node list entry 256 for application object A references node list
18 entry 262 of data object O as both a predecessor and a successor. This correlates to cycle
19 dependency in that the data object O must be flushed both before (or not later than) and
20 after (or not earlier than) application object A.

21 The cache manager recognizes this cyclic condition when it occurs, or when the
22 cache manager goes to flush the application object A. For purposes of continuing
23 discussion, suppose the cache manager decides to flush the application object A. The
24 cache manager proceeds down A's node list, which contains the single entry 256, and
25

discovers the cycle dependency. When a cycle between the intermediate write graph nodes 251 and 253 is discovered, the nodes 251 and 253 are collapsed into a single node.

Fig. 21 shows a write graph 209 having a combined write graph node 255 formed by collapsing nodes 251 and 253 into one node following the execute operation 196 (i.e., Ex_{196}). The node ID field 240' of object O's node list entry 262 is switched from "null" to reference an entry 257 in the node table 246. Additionally, the node ID field 240' of object A's node list entry 256 is changed from "null" to reference the entry 257. The node table entry 257 lists all intermediate graph nodes (identified by pairs <Object, Object State ID>) from which it is constituted via the collapse. In this example, the node table entry 257 identifies the node 251 as <A, 4> and the node 253 as <O, 2>.

To break the cycle dependency and flush the object A by itself, the cache manager first installs all write graph nodes preceding the object A. In this case, the only real predecessor node (which is a node of the intermediate write graph) contains object O, which forms the cycle dependency with A and hence is to be flushed simultaneously with the application object A. The cache manager then blindly writes the data object O listed in the predecessor list 242' of object A's node list entry 256 to the stable log. That is, the values of the data object at state 2 (i.e., O_2) are posted to the stable log. This is shown in Fig. 16 as the blind write 216, which results in a log record containing the value O_2 .

Fig. 22 shows the blind write operation 216 of data state O_3 and a corresponding write graph 211. The write graph 211 contains three nodes: a node 259 containing exposed object A and unexposed object O, a node 261 containing exposed object B, and a node 263 containing exposed object O. As a result of the write operation 216, a second entry 264 is added to the node list for data object O. This second entry 264 has a last field 238''' set to a state ID of 3, a node ID field 240''' set to null, a predecessor list field 242''' set to reference the node 261 containing application object B as a real predecessor

node, and a successor list field 244''' set to null. The node list entry 256 for object A is also updated following the write operation 216. The last field 238' has been updated to A's last state ID of 4, and the predecessor field 242' is updated to identify the node 261 containing application object B as a predecessor node.

Notice that the node ID fields in A's node list entry 256 and O's node list entry 262 remaining pointing to entry 257 in the node table 246. The cycles have not yet disappeared. The node for data object O in the cycle is no longer the last node for object O, so object O is not in vars(257). But the operations that previously wrote data object O are still in node 259, and this is what is captured by having the node IDs continue to reference 257. The blind write operation 216 rendered object O "unexposed" in node 259 and creates a new intermediate write graph node 263 for data object O.

A node list entry 266 for the application object B is also shown in Fig. 22. This entry 266 reflects the node 261 that was created by the read operation 198 in Fig. 16, prior to the blind write operation 216. The object B's node list entry 266 has a last field 238''' set to B's last state ID of 2, a node ID field 240''' set to null, a predecessor list field 242''' set to null, and a successor list field 244''' set to identify the nodes 259 and 263.

Notice that the predecessor list field 242' in object A's entry 256 still contains reference to the data object O. Predecessors are only removed when a flush occurs, and not as a result of the blind write operation 216. This is because there can be other operations on other objects that continue to depend on the prior version of the just logged object. However, the blind write does remove the blind written object O from the objects that need to be flushed simultaneously with object A.

Suppose the cache manager wishes to flush application object A. Before doing that, the node containing A must not have predecessors in the write graph. Thus, the

1 cache manager must first flush B to remove B's node 266 from the write graph. Next, the
2 cache manager flushes the application object A, thereby installing the operations 190-196
3 contained in node 259 of Fig. 22, which is represented by node table entry 257 of Figure
4 21. Fig. 23 shows the results of flushing application object A. The object O's node list
5 entry 262 which contains reference to the node 259 via node table entry 257 that it
6 references via its node ID field is removed as these states are now installed. The
7 successor list field 244' in A's entry 256 is updated to remove all successors since A has
8 now been installed. That is, the flushing of A leaves it's node list entry 256 with no
9 successors. Accordingly, this flushing operation removes the intermediate graph cycle
10 dependency as the node list entry 256 for application object A no longer contains
11 reference to data object O in either the successor or predecessor list fields.

12 The write optimization techniques described in this section are beneficial because
13 they eliminate having to always post the written values to the log. This greatly reduces
14 the processing time cost during normal operation, at the expense of more costly recovery
15 processing. With the optimization, the log only contains information to reference its
16 source object and the state ID of the values that are written. While this reduces the
17 amount of data to be logged, the write optimization techniques introduce dependencies
18 between objects, and often troubling cycle dependencies. The cache manager tracks
19 dependencies via an object table and is configured to recognize cycle dependencies.

20 When a cycle dependency is realized, the cache manager initiates a blind write of
21 one or more objects involved in the cycle to place the object's values on the stable log.
22 This step breaks the cycle. Thereafter, the cache manager flushes the objects according to
23 an acyclic flushing sequence that pays attention to any predecessor objects that first
24 require flushing. The acyclic flushing sequence is structured such that the object that is
25 removed from the cycle dependency by the blind write is flushed to the stable database

1 after the other object of the original cycle dependency. In other words, the object that is
2 not removed from the cycle dependency by the blind write is flushed to the stable
3 database before the object that is removed from the cycle dependency is flushed to the
4 stable database. If multiple blind writes are used to render multiple objects in a multi-
5 object node unexposed, thereby removing them from the multi-object node, these objects
6 that are unexposed and no longer in the original node are flushed to the stable database
7 after the exposed object(s) that remain in the original node are flushed to the stable
8 database.

9 As described, the present invention breaks up atomic flush sets, regardless of
10 whether they are produced by cyclic flush dependencies or otherwise, such as by one
11 operation writing two objects, as described above with respect to Figs. 28A – 28C, and
12 requiring that the objects be flushed atomically.

13 It should be noted that the data structures used by the cache manager as described
14 in accordance with the present invention are directed to a single updated object per
15 operation because an object table entry is used to represent, at least some of the time, a
16 write graph node. However, the current invention can work with other cache manager
17 data structures that permit operations to update more than a single object per operation.

18 The object table 222 of Fig. 18, similar to that described with respect to Fig. 11, is
19 used to manage the acyclic flushing sequence of the various objects. The objects can be
20 application and/or data objects. As described above, each entry 224 in the object table
21 222 has a node field 234 that contains an index to a separate node list 236 of intermediate
22 graph nodes. Each entry in the node list 236 has a predecessor list 242 and a successor
23 list 244. These lists are used to track the flushing sequences of the various nodes; i.e.,
24 these lists determine which nodes and their object(s) should be flushed before the subject
25

node and object(s). When an object is flushed, the object is removed from its successors' predecessor list entries.

Recovery Optimization

During recovery, the database computer system can invoke a conventional recovery manager to recover the application state and object state at the instance of the crash. The conventional recovery manager retrieves the most recently flushed data objects and application objects in the stable database. The recovery manager then replays the stable log, beginning at a point known to be earlier than the oldest logged operation that was not yet installed. For this conventional physiological operation recovery, the recovery manager compares the state ID of each logged operation in the stable log with the state ID of a retrieved data object or application object. If the state ID of the logged operation is later than the state ID of the stable object, the recovery manager redoes that logged operation.

Fig. 24 pertains to a conventional recovery approach that can be used in conjunction with aspects of this invention. Fig. 24 shows an excerpt from a stable log, referenced generally as number 270, having a series of log records posted as a result of computer application operations. For purposes of discussion, assume that the log records in log record 270 pertain only to data object O and application object A. Only log records for data object O are described.

The log excerpt shows five log records 272-280 pertaining to operations that affect data object O. The first log record 272 contains the object ID "O" and state ID "n" to reflect that the data object O was written or updated to a state tagged with a state ID of "n." Two subsequent log record 274 and 276 reflect that the data object O is written two more times, at states n+g and n+h. A fourth log record 278 reflects that the entire value

1 for the data object O at state $n+h$ (i.e., O_{n+h}) is written to the stable log, as is the case for a
2 blind write operation, at a state ID of “ $n+i$ ”.

3 Each log record is assigned a log sequence number (LSN). The LSN is a
4 monotonically increasing number that is tagged to each log record to identify the order in
5 which the log records are created in the log. Typically, the LSN is used as the state ID,
6 making the state ID and LSN the same. The LSN for the log records 272-278, for
7 instance, are n , $n+g$, $n+h$, and $n+i$.

8 Suppose that the cache manager flushes the data object at its state “ n ” (i.e., O_n) to
9 the non-volatile database. This event is recorded as log record 280 that identifies the data
10 object O as having been flushed. All log records for the data object O that precede log
11 record 272 are no longer needed for replay during recovery. In fact, log record 272 is not
12 really needed for replay because it simply identifies the exact object state that is present
13 in the database. Rather, the first meaningful log record for recovery purposes is the first
14 log record reflecting an operation that updates the data object O, thereby changing its
15 state, without the updated data object O being flushed to install the operation. In this
16 example, the first meaningful log record is record 274.

17 At the time that data object O is flushed, the cache manager marks object O as
18 clean (the dirty flag is reset) in the cache. When O is updated at log record 274, the cache
19 manager sets a recovery log sequence number (rLSN) to identify the log record 274 as the
20 starting point for replay of object O during recovery.

21 Each object has its own rLSN. In this example, data object O has an rLSN and
22 application object A has a separate rLSN (not shown). During recovery, the recovery
23 manager examines the last checkpoint record on the stable log, which contains initial
24 values of rLSNs for all dirty objects as of the time of the checkpoint. Subsequent logging
25 of flushes merely updates which objects are clean or dirty and advances rLSNs as these

changes occur. Alternatively, the checkpoint record can indicate the value of the minimum rLSN, so that the individual rLSNs can be recomputed based on the updates to objects and their flushing. But in this case, it needs to at least bound the rLSN before proceeding. The recovery manager then begins its redo test at the minimum rLSN_{min}. The recovery manager examines every record thereafter to determine whether to replay the operation. This portion of the log after the rLSN_{min} is known as the “active log tail.”

A shortcoming of this conventional recovery technique is that the recovery manager can end up replaying many operations that are unnecessary for recovery. As an example, the lifetimes of some application objects and data objects tend to be short and once terminated or deleted the objects no longer need recovery. If a system failure occurs after an object has terminated, but while that object's updates remain on the active log tail, the recovery manager still redoes the operations for that object starting from the last stable version of the object. If the object's state was never written to stable memory, all updates reflected in the log records are redone. Unfortunately, the replayed operations for these terminated or deleted objects are unnecessary, and can add substantially to recovery time.

Accordingly, an aspect of this invention is to optimize recovery to avoid replaying operations that are rendered obsolete by subsequent operations. In general, the recovery optimization technique involves advancing an object's rLSN to a log record later in the stable log that reflects the object at a state in which the operations that have written that object state are installed. Normally, flushing the object to non-volatile memory is what installs earlier operations and so capturing the change in rLSN could be done by logging the object flushes. But when dealing with objects that are “unexposed” in the write graph, the operations leading to a particular object can be installed without that object itself being flushed.

Recall the discussion from Figs. 16 and 17. A blind write operation posted the value of data object O to the stable log and thereby rendered the data object O “unexposed” in the write graph node containing application object A, meaning that the prior value of O was no longer needed at the time when the write graph node is installed. The flushing of application object A installed all operations (i.e., R_{190} , Ex_{192} , W_{194} , and Ex_{196}), including the write operation W_{194} that had written the data object O, even though the data object O itself was not flushed. Preferably, however, the data object O is flushed to the stable database to make cache management effective by providing a place in the non-volatile database from which the object value can be retrieved should its value be dropped from the cache.

Fig. 25 shows an example of the recovery optimization technique for the same stable log 270. Suppose that log record 278 represents a blind write operation in which the cache manager posts the values of data object O at state ID of “n+h” (i.e., O_{n+h}) to the stable log. The blind write renders the data object O “unexposed” in the write graph node containing both objects A and O, as described above with respect to Fig. 17.

Sometime after the blind write operation, the cache manager flushes the “exposed” application object A at state “m” (i.e., A_m) to install all operations in the write graph node, including any operations that have written the data object O. The blind write and subsequent flushing of application object A renders all operations that wrote the “unexposed” data object O as part of the operations associated with the node for application object A (e.g., log records 274 and 276) unnecessary for recovery.

The cache manager advances the $rLSN_A$ for the “exposed” application object A (not shown in this figure) because all preceding operations affecting A are now installed, akin to the customary case shown in Fig. 24. Similarly, the cache manager advances the $rLSN_O$ for the “unexposed” data object O from its original point at log record 274 to the

new location after log record 278 as if the unexposed data object O had also been flushed. Record 278 is the next log record, after the records for the installed operations, that contains an operation writing data object O. The rLSN of object O is logged as a record 284 to reference the log record 278 with the log sequence number of $n+i$. In this manner, the recovery manger treats “unexposed” objects as if they had been flushed as of their last update in the write graph node being installed by the flushing of the node’s exposed variable(s). By logging the rLSN, recovery for O can begin at log record 278.

The rLSN is recorded in the cache manager’s object table. Fig. 26 shows a cache manager 290 and object table 292 that are similar in structure to that shown in Fig. 18. However, in Fig. 26, an entry 294 for data object O is modified to include an rLSN field 296, which identifies the LSN of the next log record that contains an operation writing data object O (in this case, $n+i$ for log record 278). This log record contains the first update to O since it was installed. The dirty flag remains set to indicate that the data object has been updated since its last value was installed.

To ensure that the object table is recoverable, and hence the rLSNs, the rLSN is also posted to the stable log as its own log record. Fig. 25 shows a log record 284 that contains identification of the rLSN for object O.

During recovery, the recovery manager 71 performs two passes: (1) an analysis pass and (2) a redo pass. During the analysis pass, the recovery manager scans the active log tail to locate the rLSNs for all objects. In this example, the $rLSN_O$ for data object O references an LSN of $n+i$ for log record 278. The recovery manager next identifies the minimum recovery log sequence number $rLSN_{min}$, similar to the conventional method described above. However, because the rLSNs have been advanced using the recovery optimization techniques, the $rLSN_{min}$ could be much later in the log as compared to the

1 conventional recovery method, thereby avoiding the replay of operations that are
2 unnecessary for recovery.

3 During the redo pass, the recovery manager examines all operations on the log
4 beginning at the $rLSN_{min}$. More particularly, the recovery manager performs the
5 following redo test for each log record in the stable log that follows $rLSN_{min}$:

- 6
7 1. If the LSN of the log record of object O is less than the $rLSN_O$ for object O
8 (meaning that the operation referenced by that record occurred before the log
9 record tagged with $rLSN_O$), the redo test is false and the operation in the log
10 record is not replayed. This condition indicates that the operation is installed
11 and the object is not exposed.
12
- 13 2. If the LSN of the log record is greater than or equal to $rLSN_O$ (meaning that it
14 occurred after the last logged installation of object O), the redo test may be
15 true. Data object O is read from stable storage and the LSN stored with O is
16 then used as $rLSN$. The redo test is then performed using the new $rLSN$, and
17 if true, the operation in the log record is replayed. This condition indicates
18 that the operation is not installed and the variable is exposed.
19

20 The redo pass rebuilds the object table, complete with $rLSNs$ for each object
21 during the analysis phase. So long as the LSN of the log record for an operation
22 involving writing O is less than object O's $rLSN_O$, the redo test returns false and the
23 operation is ignored.

24 Once the log record for an object O is greater than or equal to its $rLSN_O$ (as seen
25 in the recovered object table), the stable version of object O (if there is one) is read to

compare the log LSN with the LSN stored with the value of O. (This can be higher than the $rLSN_O$ should the system have failed between the time data object O was last flushed and the time the change to its $rLSN$ resulting from that flush was posted to the stable log. The $rLSN_O$ is set to the stable LSN of the value of O when this occurs.

One situation where the recovery optimization technique is helpful concerns short-lived applications that initiate, then execute and write their results, and terminate. Fig. 27 shows an exemplary excerpt from a stable log 300 having log records 302-308 for the short-lived application. The log records 302-308 correspond to the four operations: initiate, execute, write, and terminate.

The short-lived applications do not need to be replayed during recovery (assuming the results written by the application are logged or contained in a stable object). Accordingly, for such short-lived applications, the cache manager posts the $rLSN_A$ for the application object A to the last operation for object A, i.e., the terminate operation recorded in log record 308. The $rLSN_A$ is posted to the stable log as record 310. Note that the $rLSN$ cannot be advanced simply because of the terminate operation 308, as versions of A may still be needed, e.g. to recover object O. During the redo pass of recovery, the recovery manager proceeds to the $rLSN_A$ for that application object and finds that the log record pertains to a terminate operation, which does not need to be redone.

As a result, the recovery manager avoids replaying the set of operations for the short-lived application object A. When application A has written an object O, if the value of O that A wrote has been installed (whether by explicit flush or because it is no longer exposed), A does not need to be recovered so that O can be recovered. Further, if application A reads data object O, but application A has been installed, either by flushing or because A's state is no longer exposed (e.g. it might be terminated or it might have

1 been written to the log), then object O need not be recovered so as to recover application
2 object A. The fact that application object A terminated is not sufficient to dispense with
3 recovering object A as it may be needed to reconstruct objects that it wrote. However,
4 when the terminate operation for A is installed (and at that point, A is not exposed), then
5 we advance A's rLSN to indicate that A's recovery is no longer needed.

6 It should also be noted that rLSN's can be advanced without actually writing them
7 to the log, though logging them in this way greatly simplifies the analysis pass of
8 recovery. Without logging rLSN's, but continuing to log the flushing of objects, the
9 analysis pass must examine each logged operation and re-create the write graph for the
10 operations as they are encountered, based on the objects read and written by each
11 operation. This permits the analysis pass of recovery to determine when the flushing of a
12 variable installs unexposed objects as well. That permits it to advance the rLSN's for
13 these objects.

14 The invention has been described in language more or less specific as to structure
15 and method features. It is to be understood, however, that the invention is not limited to
16 the specific features described, since the means herein disclosed comprise exemplary
17 forms of putting the invention into effect. The invention is, therefore, claimed in any of
18 its forms or modifications within the proper scope of the appended claims appropriately
19 interpreted in accordance with the doctrine of equivalents and other applicable judicial
20 doctrines.

CLAIMS

1. In a database computer system having a non-volatile memory, a volatile main memory, and a first object which executes from the main memory, wherein the non-volatile memory includes a stable log, a computer-implemented method comprising the following steps:

executing the first object to perform operations which read data from, and write data to, a second object;

posting to the stable log a log record for each operation involving the reading or writing of data, the log record containing a reference to either the first object or the second object to identify that referenced object as a source for the data that is read from or written to;

establishing flush order dependencies between the first object and the second object, wherein some of the flush order dependencies become cyclic indicating a condition in which the first object should be flushed not later than the second object and the second object should be flushed not later than the first object;

detecting a dependency cycle;

following detection of the dependency cycle, writing one of the first object and the second object to the stable log to break the dependency cycle;

flushing the other of the first object and the second object to the non-volatile memory; and

flushing the object written to the stable log to the non-volatile memory.

1 2. A computer-implemented method as recited in claim 1, wherein the first
2 object is an application object and the second object is a data object.

3
4 3. A computer-implemented method as recited in claim 2, wherein the writing
5 step writes the data object to the stable log to break the dependency cycle, and the
6 flushing steps flush the application object to the non-volatile memory prior to flushing
7 the data object to the non-volatile memory.

8
9 4. A computer-implemented method as recited in claim 1, wherein the writing
10 step forms a flush dependency edge between the first object and the second object.

11
12 5. A computer programmed to perform the steps of the computer-
13 implemented method as recited in claim 1.

14
15 6. A computer-readable memory that directs a computer to perform the steps
16 in the method as recited in claim 1.

17
18 7. In a database computer system having a cache manager which occasionally
19 flushes objects from a volatile main memory to a non-volatile memory to preserve those
20 objects in the event of a system crash, and wherein a dependency cycle exists between at
21 least two objects such that the two objects should be flushed simultaneously, a computer-
22 implemented method comprising the following steps:

23 detecting a dependency cycle;

24 writing one of the two objects to the stable log to break the dependency cycle;

25 flushing the other of the two objects to the non-volatile memory; and

1 flushing the object that has been written to the stable log to the non-volatile
2 memory.

3
4 **8.** A computer-implemented method as recited in claim 7, wherein one of the
5 two objects is an application object and the other of the two objects is a data object.

6
7 **9.** A computer-implemented method as recited in claim 8, wherein the writing
8 step writes the data object to the stable log to break the dependency cycle, and the
9 flushing steps flush the application object to the non-volatile memory prior to flushing
10 the data object to the non-volatile memory.

11
12 **10.** A computer-implemented method as recited in claim 7, wherein the
13 writing step establishes a flush dependency edge between the two objects.

14
15 **11.** A computer programmed to perform the steps of the computer-
16 implemented method as recited in claim 7.

17
18 **12.** A computer-readable memory that directs a computer to perform the steps
19 in the method as recited in claim 7.

20
21 **13.** A database computer system comprising:
22 a volatile main memory;
23 a non-volatile memory that persists across a system crash;
24 a processing unit coupled to the main memory and the non-volatile memory;
25

1 a first object stored in the volatile main memory and executable on the processing
2 unit;

3 a resource manager which interacts with the first object to mediate
4 communication between the first object and a second object so that, during an operation,
5 the resource manager writes data from the first object to the second object;

6 the resource manager being configured to log, in a log record on the non-volatile
7 memory, a reference to the first object to identify the first object as a source for the data
8 that was written to the second object; and

9 the resource manager including a cache manager for establishing a flush order
10 dependency between the first object and the second object as a result of the operation and
11 managing a flushing order in which the first object and the second object are occasionally
12 flushed to the non-volatile memory according to the flush order dependency,

13 wherein the operation results in a dependency cycle between the first object and
14 the second object indicating that the first and second objects should be flushed
15 simultaneously, the cache manager being configured to detect the cycle dependency and
16 in response to the detection, to write one of the first object or the second object as a log
17 record to the non-volatile memory to break the dependency cycle so that the first object
18 and second object can be flushed to the non-volatile memory in a sequential manner, to
19 flush the other of the first object and the second object to the non-volatile memory, and
20 then to flush the one of the first object or the second object to the non-volatile memory.

21
22 **14.** A database system as recited in claim 13, wherein the first object is an
23 application object and the second object is a data object.
24
25

1 **15.** A database system as recited in claim 14, wherein the cache manager is
2 configured to write the data object to the stable log to break the dependency cycle, and to
3 flush the application object to the non-volatile memory and then to flush the data object
4 to the non-volatile memory.

5
6 **16.** A database system as recited in claim 13, wherein the cache manager is
7 configured to establish a flush dependency edge between the first object and the second
8 object to break the dependency cycle.

9
10 **17.** For use on a database computer system having a non-volatile memory, a
11 volatile main memory, a processing unit, a first object stored in the main memory and
12 executed on the processing unit, and a second object stored in the main memory, a cache
13 manager executable on the processor to manage flushing of the first object and the second
14 object from the main memory to the non-volatile memory, the cache manager being
15 configured to detect any cycle dependency between the first object and the second object
16 indicating that the first and second objects should be flushed simultaneously, wherein in
17 response to detecting a cycle dependency, the cache manager writes one of the first object
18 and the second object as a log record to the non-volatile memory to break the dependency
19 cycle so that the first object and second object can be flushed to the non-volatile memory
20 in a sequential manner, flushes the other of the first object and the second object to the
21 non-volatile memory, and then flushes the one of the first object or the second object to
22 the non-volatile memory.

1 **18.** A cache manager as recited in claim 17, wherein the first object is an
2 application object and the second object is a data object.

3
4 **19.** A cache manager as recited in claim 18, wherein the data object is written
5 as a log record to break the dependency cycle, and the application object is flushed to the
6 non-volatile memory prior to flushing the data object to the non-volatile memory.

7
8 **20.** A cache manager as recited in claim 17, wherein a flush dependency edge
9 between the first object and the second object is established when the dependency cycle is
10 broken.

11
12 **21.** In a database computer system having a non-volatile memory, a volatile
13 main memory, and a first object which executes from the main memory, wherein the non-
14 volatile memory includes a stable log, a computer-implemented method comprising the
15 following steps:

16 executing the first object to perform operations which read data from, and write
17 data to, a second object;

18 posting to the stable log a log record for each operation involving the reading or
19 writing of data, the log record containing a reference to either the first object or the
20 second object to identify that referenced object as a source for the data that is read from
21 or written to;

22 detecting an atomic flush set comprising the first object and the second object,
23 wherein the atomic flush set indicates a condition in which the first object should be
24 flushed not later than the second object and the second object should be flushed not later
25 than the first object;

1 following detection of the atomic flush set, writing one of the first object and the
2 second object to the stable log to break up the atomic flush set;
3 flushing the other of the first object and the second object to the non-volatile
4 memory; and
5 flushing the object written to the stable log to the non-volatile memory.
6

7 **22.** A computer-implemented method as recited in claim 21, wherein the first
8 object is an application object and the second object is a data object.
9

10 **23.** A computer-implemented method as recited in claim 22, wherein the
11 writing step writes the data object to the stable log to break up the atomic flush set, and
12 the flushing steps flush the application object to the non-volatile memory prior to
13 flushing the data object to the non-volatile memory.
14

15 **24.** A computer-implemented method as recited in claim 21, wherein the
16 writing step forms a flush dependency edge between the first object and the second
17 object.
18

19 **25.** A computer programmed to perform the steps of the computer-
20 implemented method as recited in claim 21.
21

22 **26.** A computer-readable memory that directs a computer to perform the steps
23 in the method as recited in claim 21.
24
25

1 **27.** In a database computer system having a cache manager which
2 occasionally flushes objects from a volatile main memory to a non-volatile memory to
3 preserve those objects in the event of a system crash, and wherein an atomic flush set
4 comprises at least two objects such that the two objects should be flushed simultaneously,
5 a computer-implemented method comprising the following steps:

6 detecting the atomic flush set;
7 writing one of the two objects to the stable log to break up the atomic flush
8 set;
9 flushing the other of the two objects to the non-volatile memory; and
10 flushing the object that has been written to the stable log to the non-
11 volatile memory.

12
13 **28.** A computer-implemented method as recited in claim 27, wherein one of
14 the two objects is an application object and the other of the two objects is a data object.

15
16 **29.** A computer-implemented method as recited in claim 28, wherein the
17 writing step writes the data object to the stable log to break up the atomic flush set, and
18 the flushing steps flush the application object to the non-volatile memory prior to
19 flushing the data object to the non-volatile memory.

20
21 **30.** A computer-implemented method as recited in claim 27, wherein the
22 writing step establishes a flush dependency edge between the two objects.

1 **31.** A computer programmed to perform the steps of the computer-
2 implemented method as recited in claim 27.

3
4 **32.** A computer-readable memory that directs a computer to perform the steps
5 in the method as recited in claim 27.

6
7 **33.** A database computer system comprising:
8 a volatile main memory;
9 a non-volatile memory that persists across a system crash;
10 a processing unit coupled to the main memory and the non-volatile memory;
11 a first object stored in the volatile main memory and executable on the processing
12 unit;

13 a resource manager which interacts with the first object to mediate
14 communication between the first object and a second object so that, during an operation,
15 the resource manager writes data from the first object to the second object;

16 the resource manager being configured to log, in a log record on the non-volatile
17 memory, a reference to the first object to identify the first object as a source for the data
18 that was written to the second object; and

19 the resource manager including a cache manager for establishing a flush order
20 dependency between the first object and the second object as a result of the operation and
21 managing a flushing order in which the first object and the second object are occasionally
22 flushed to the non-volatile memory according to the flush order dependency,

23 wherein the operation results in an atomic flush set comprising the first object and
24 the second object, the cache manager being configured to detect the atomic flush set and
25 in response to the detection, to write one of the first object or the second object as a log

1 record to the non-volatile memory to break up the atomic flush set so that the first object
2 and second object can be flushed to the non-volatile memory in a sequential manner, to
3 flush the other of the first object and the second object to the non-volatile memory, and
4 then to flush the one of the first object or the second object to the non-volatile memory.
5

6 **34.** A database system as recited in claim 33, wherein the first object is an
7 application object and the second object is a data object.
8

9 **35.** A database system as recited in claim 34, wherein the cache manager is
10 configured to write the data object to the stable log to break up the atomic flush set, and
11 to flush the application object to the non-volatile memory and then to flush the data
12 object to the non-volatile memory.
13

14 **36.** A database system as recited in claim 33, wherein the cache manager is
15 configured to establish a flush dependency edge between the first object and the second
16 object to break up the atomic flush set.
17
18
19
20
21
22
23
24
25

1 **37.** For use on a database computer system having a non-volatile memory, a
2 volatile main memory, a processing unit, a first object stored in the main memory and
3 executed on the processing unit, and a second object stored in the main memory, a cache
4 manager executable on the processor to manage flushing of the first object and the second
5 object from the main memory to the non-volatile memory, the cache manager being
6 configured to detect an atomic flush set comprising the first object and the second object,
7 wherein in response to detecting the atomic flush set, the cache manager writes one of the
8 first object and the second object as a log record to the non-volatile memory to break up
9 the atomic flush set so that the first object and second object can be flushed to the non-
10 volatile memory in a sequential manner, flushes the other of the first object and the
11 second object to the non-volatile memory, and then flushes the one of the first object or
12 the second object to the non-volatile memory.

13
14 **38.** A cache manager as recited in claim 37, wherein the first object is an
15 application object and the second object is a data object.

16
17 **39.** A cache manager as recited in claim 38, wherein the data object is written
18 as a log record to break up the atomic flush set, and the application object is flushed to
19 the non-volatile memory prior to flushing the data object to the non-volatile memory.

20
21 **40.** A cache manager as recited in claim 37, wherein a flush dependency edge
22 between the first object and the second object is established when the atomic flush set is
23 broken up.

1 **ABSTRACT**

2 This invention concerns a database computer system and method for making
3 applications recoverable from system crashes. The application state (i.e., address space)
4 is treated as a single object which can be atomically flushed in a manner akin to flushing
5 individual pages in database recovery techniques. To enable this monolithic treatment of
6 the application, executions performed by the application are mapped to logical loggable
7 operations that can be posted to the stable log. Any modifications to the application state
8 are accumulated and the application state is periodically flushed to stable storage using an
9 atomic procedure. The application recovery integrates with database recovery, and
10 effectively eliminates or at least substantially reduces the need for check pointing
11 applications. In addition, optimization techniques are described to make the read, write,
12 and recovery phases more efficient. Atomic flush sets, whether generated from cyclic
13 flush dependencies or otherwise, can be broken apart. This enables an ordered flushing
14 sequence of first flushing a first object and then flushing a second object, rather than
15 having to flush both the first and second objects simultaneously and atomically.

16
17
18
19
20
21
22
23
24
25

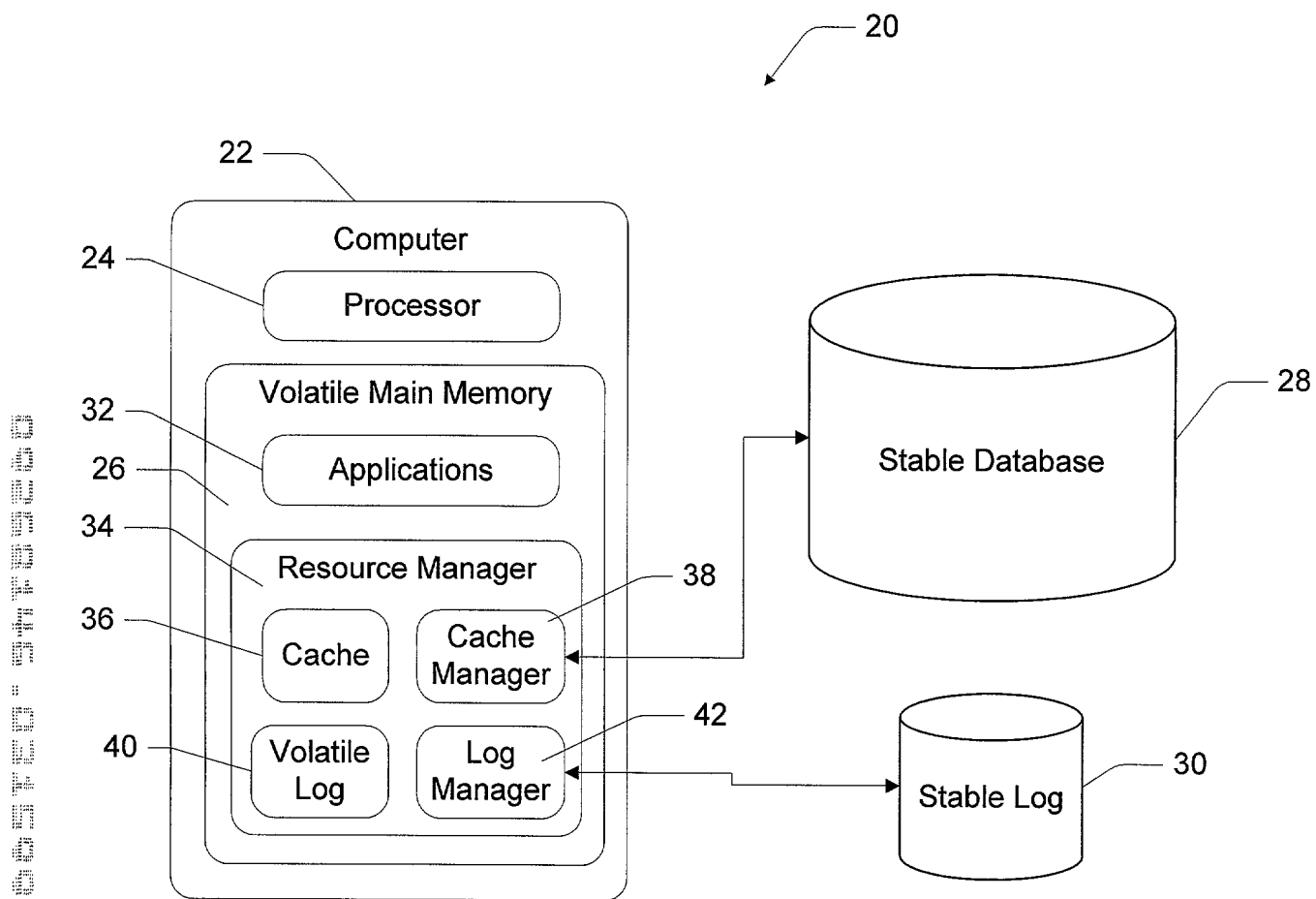


Fig. 1
Prior Art

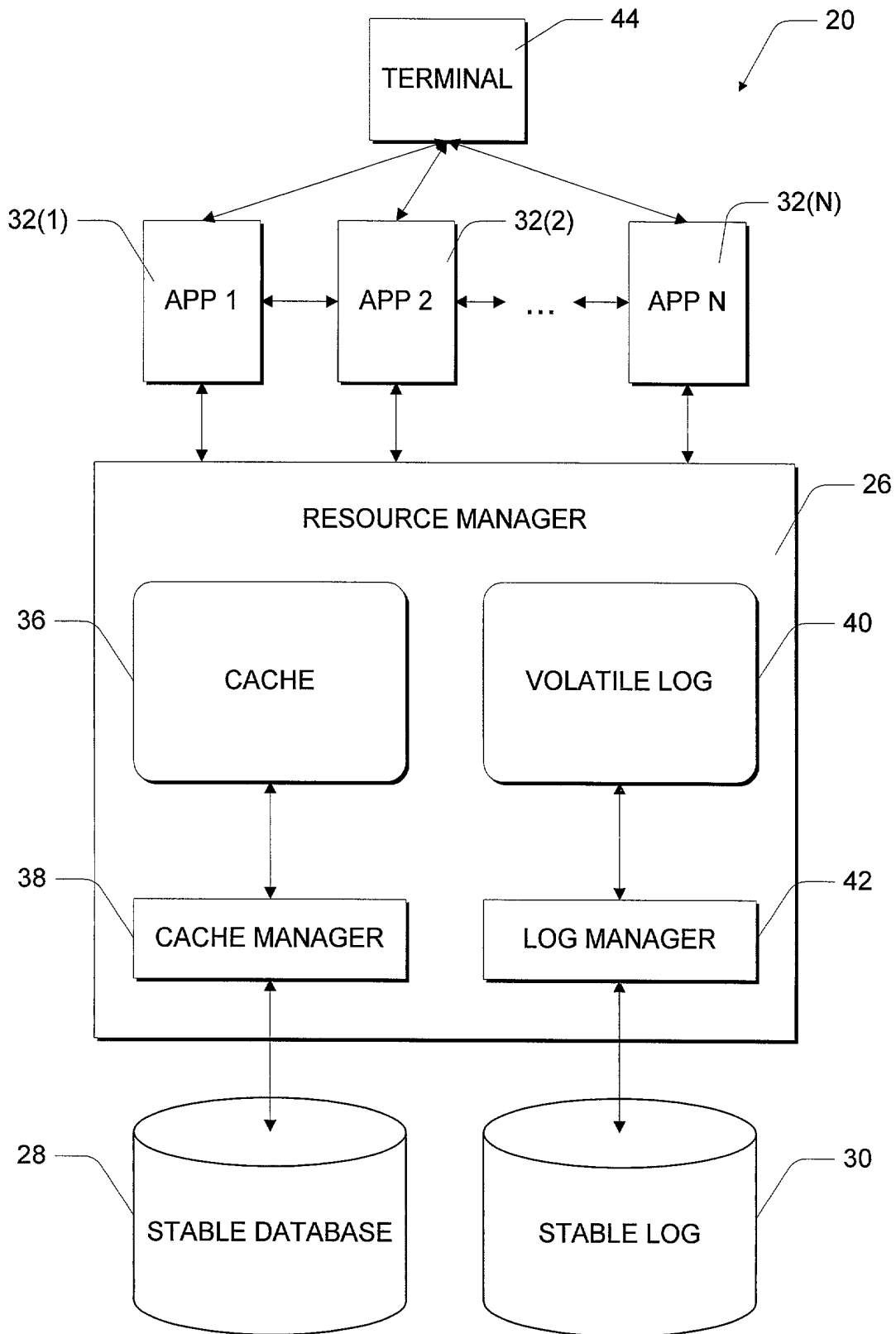


Fig. 2
Prior Art

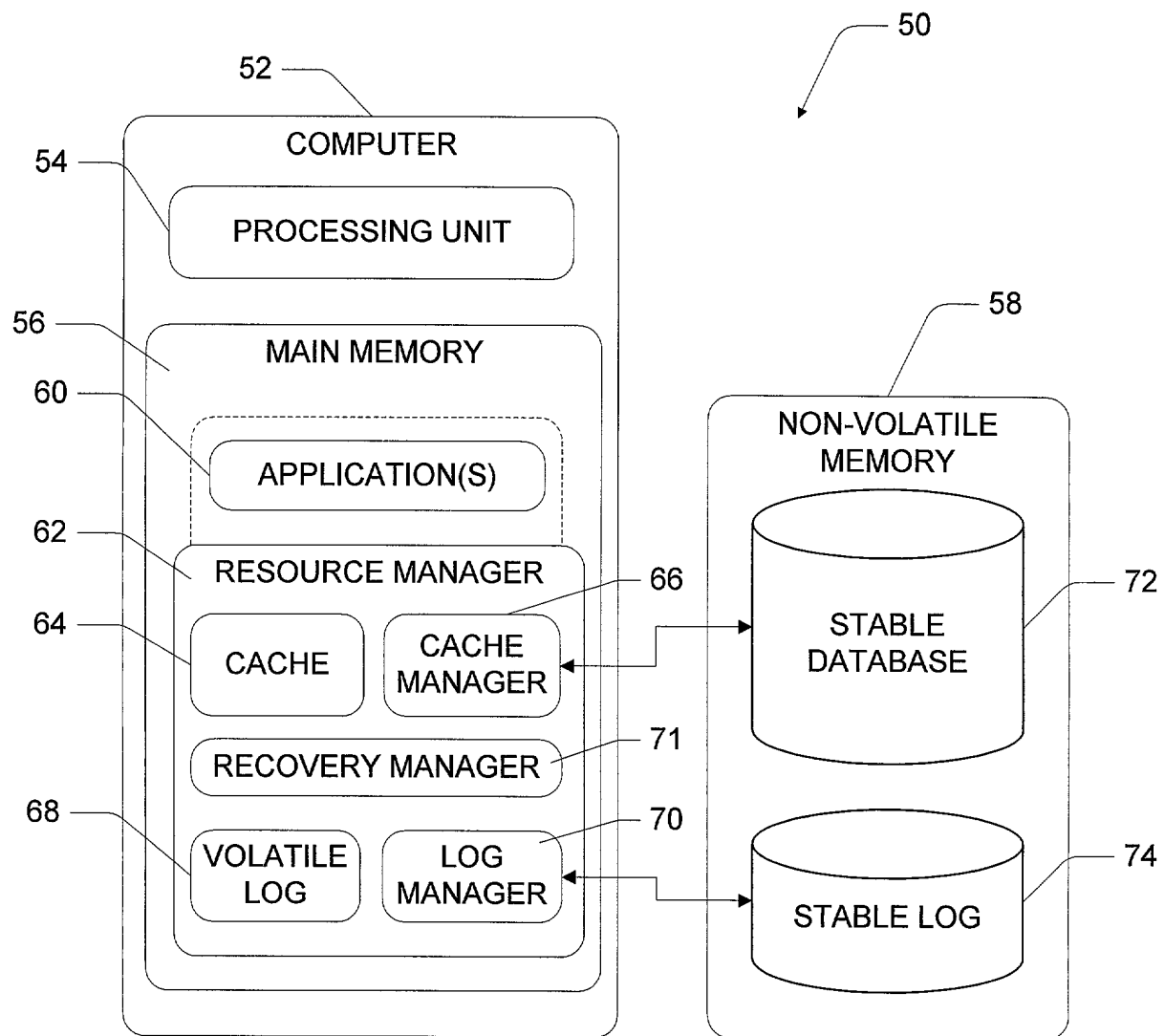


Fig. 3

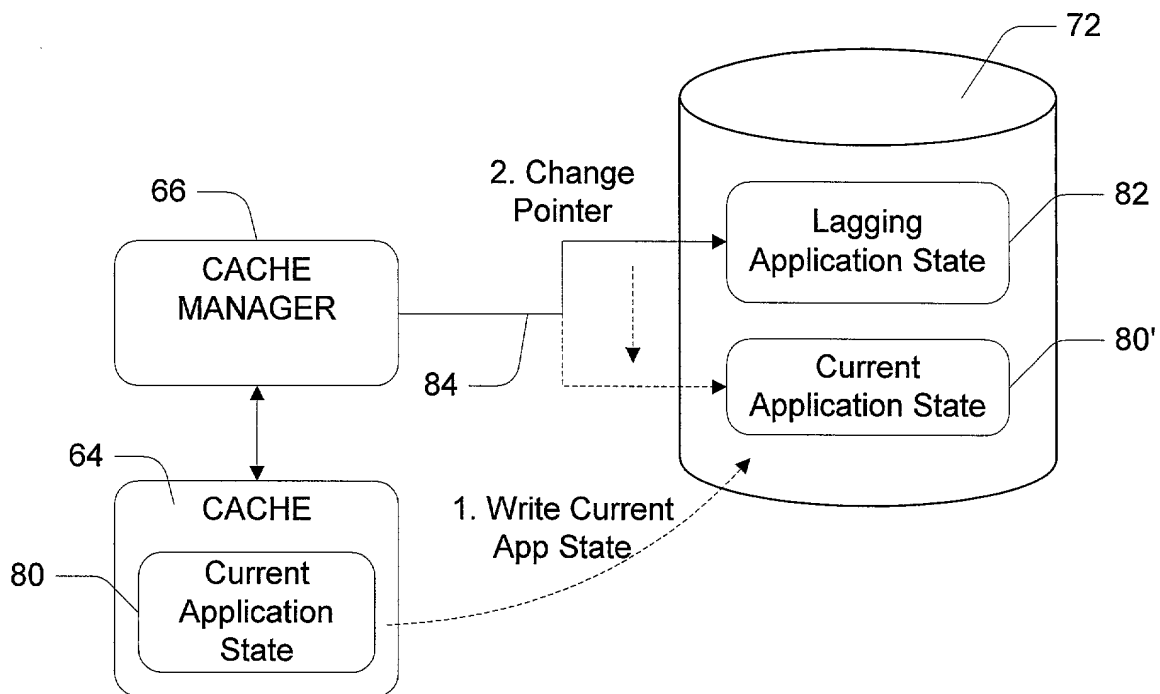


Fig. 4

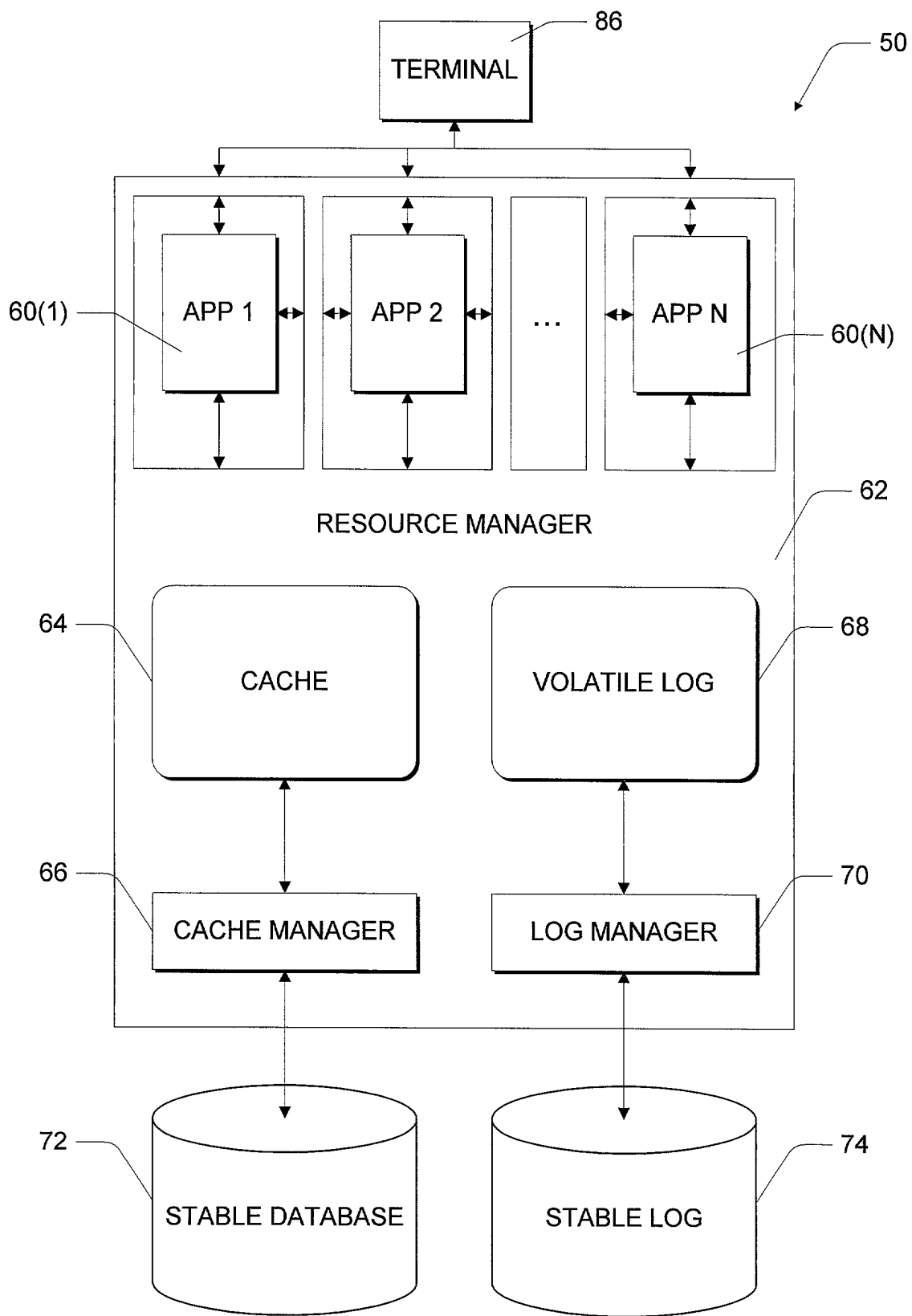


Fig. 5

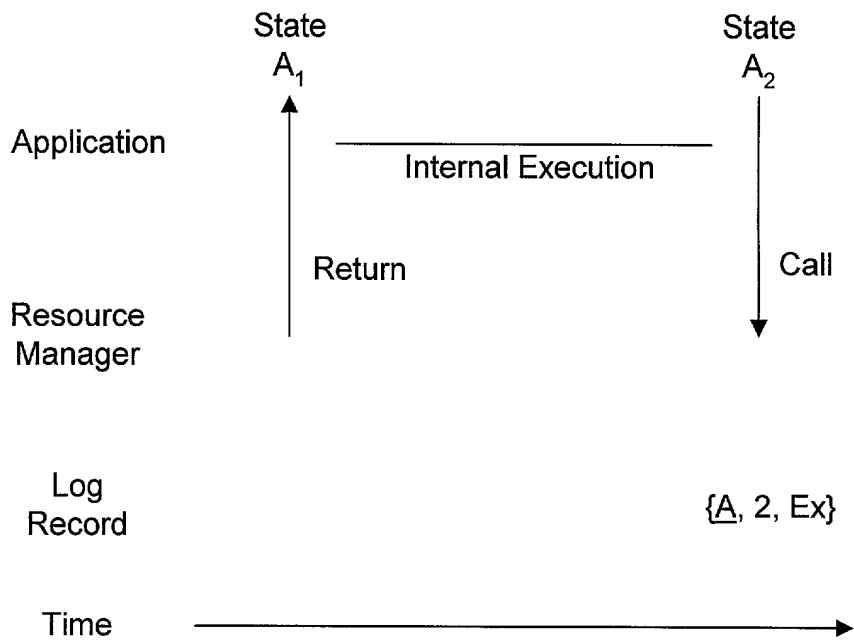


Fig. 6

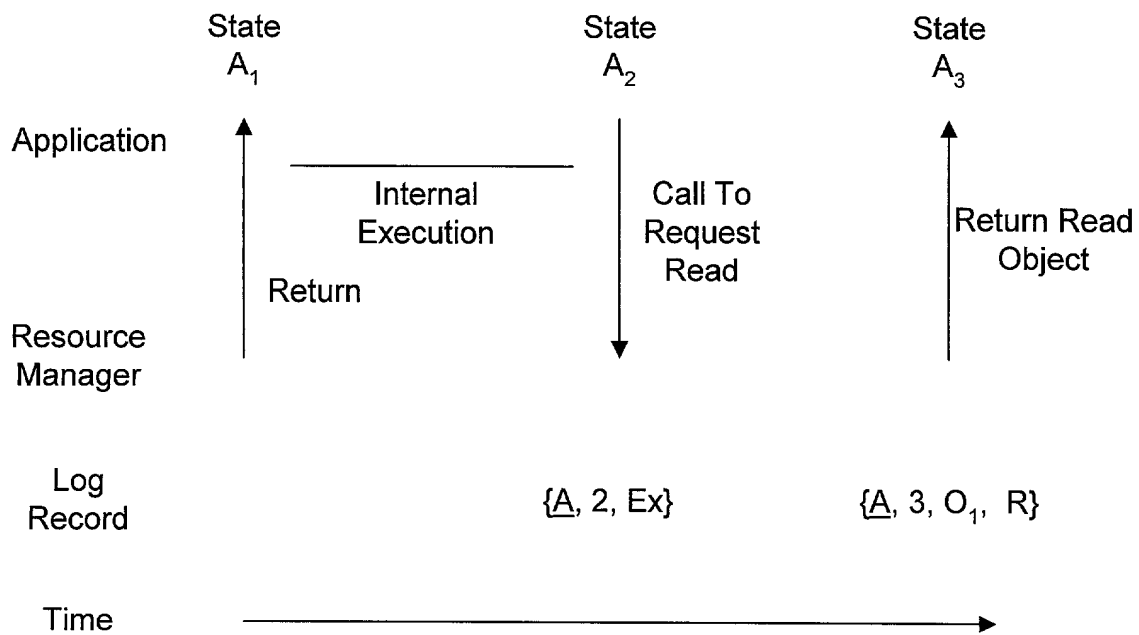


Fig. 7

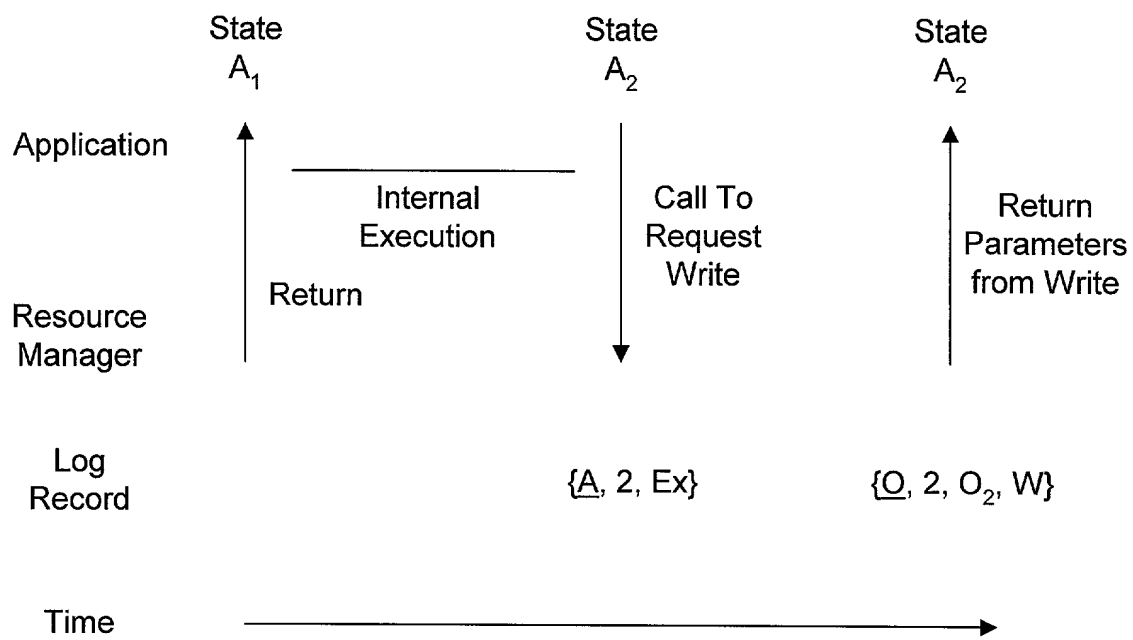


Fig. 8

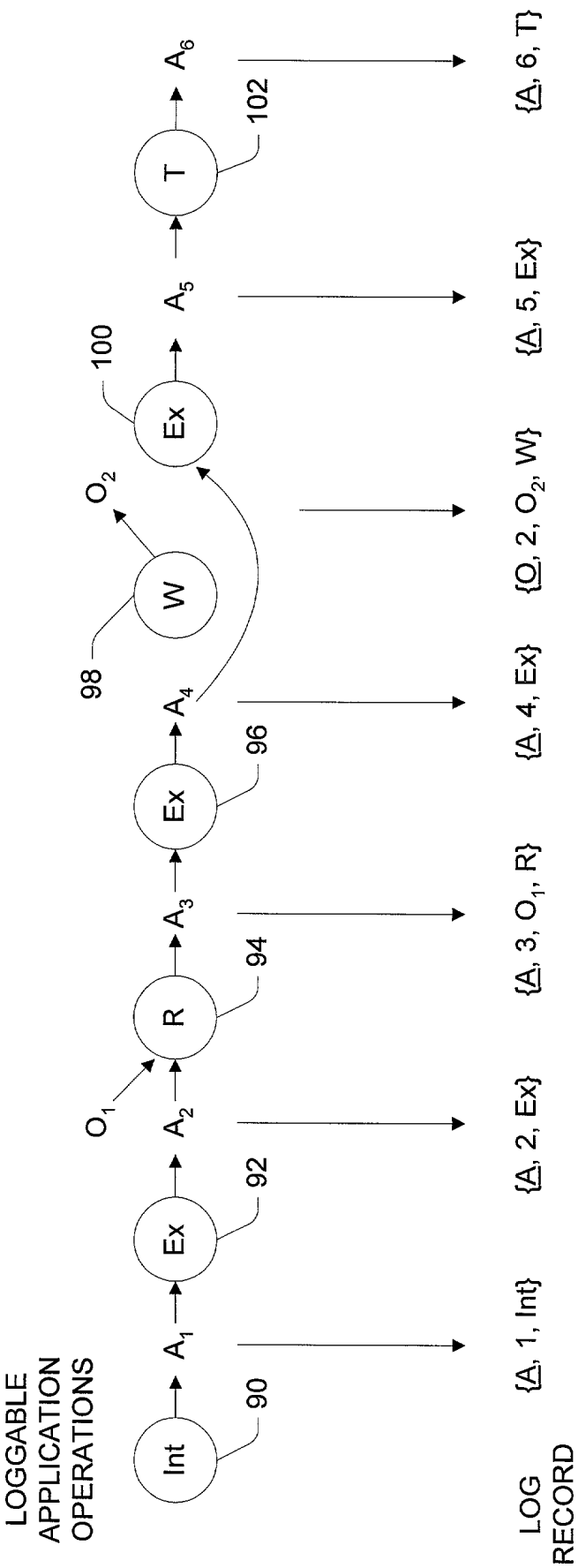


Fig. 9

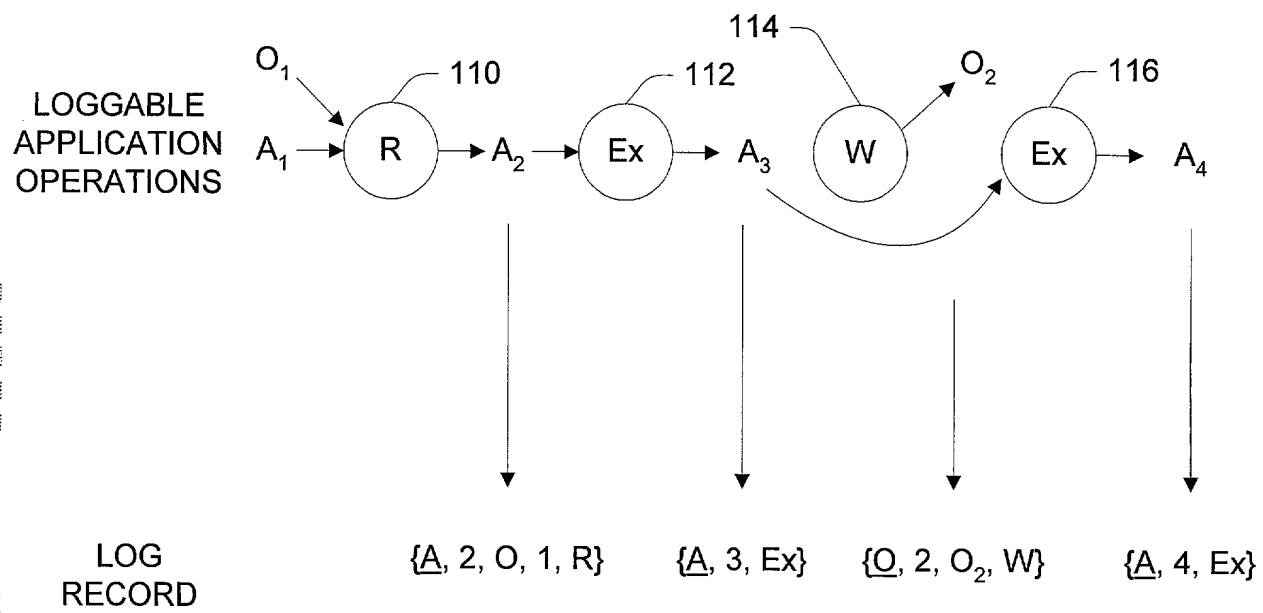


Fig. 10

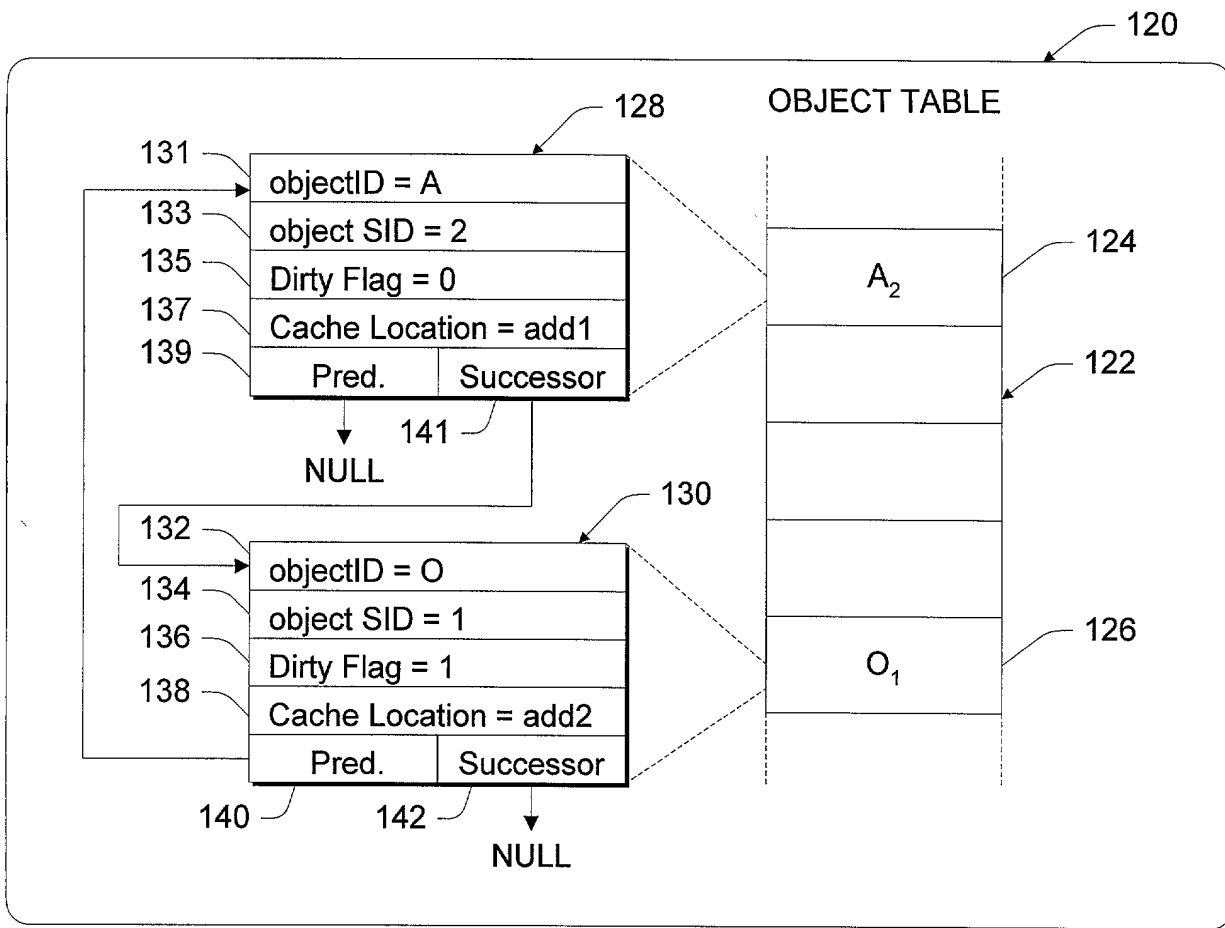


Fig. 11

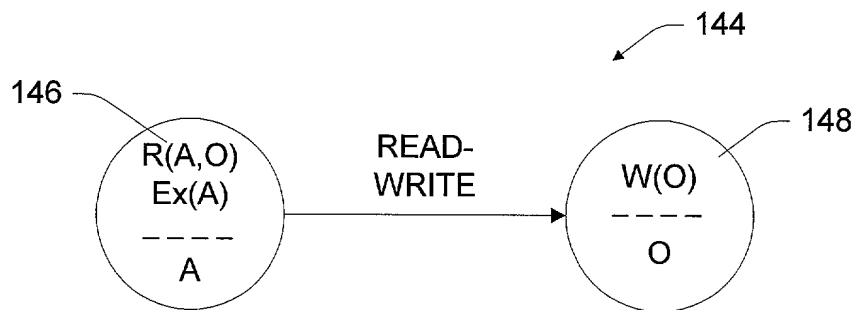


Fig. 12

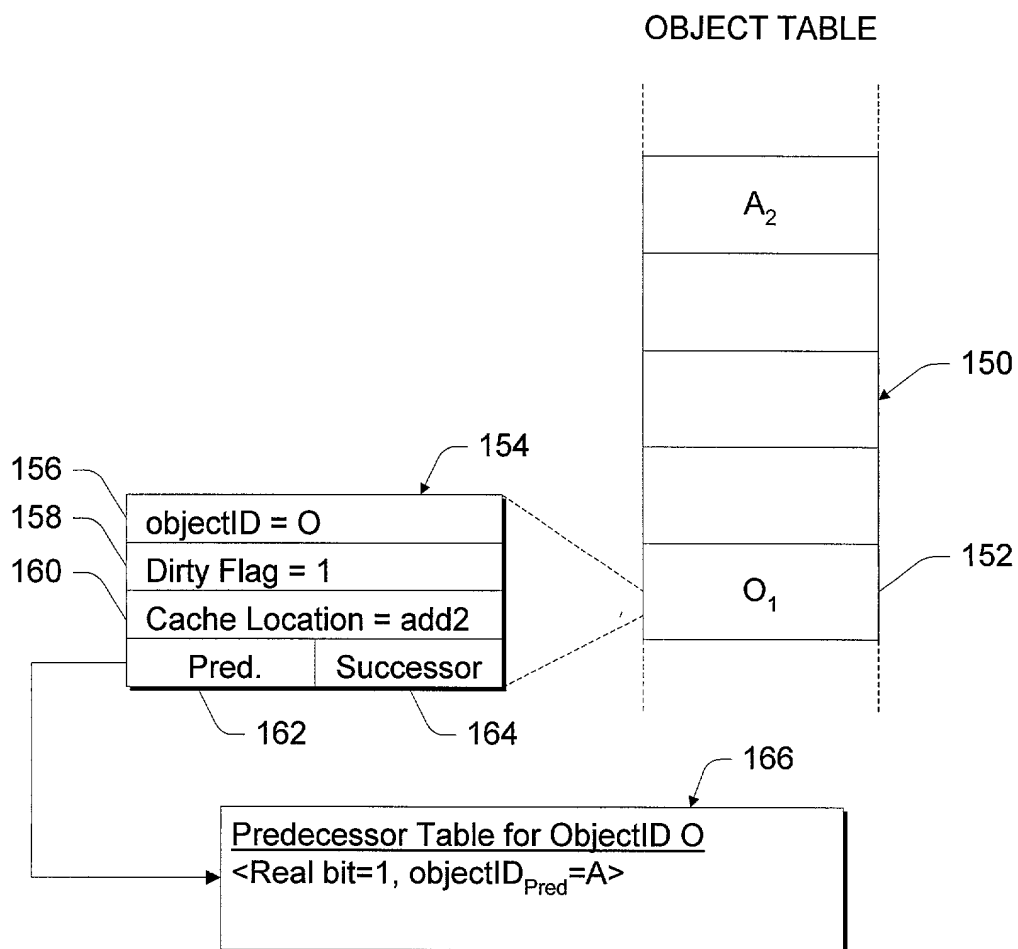


Fig. 13

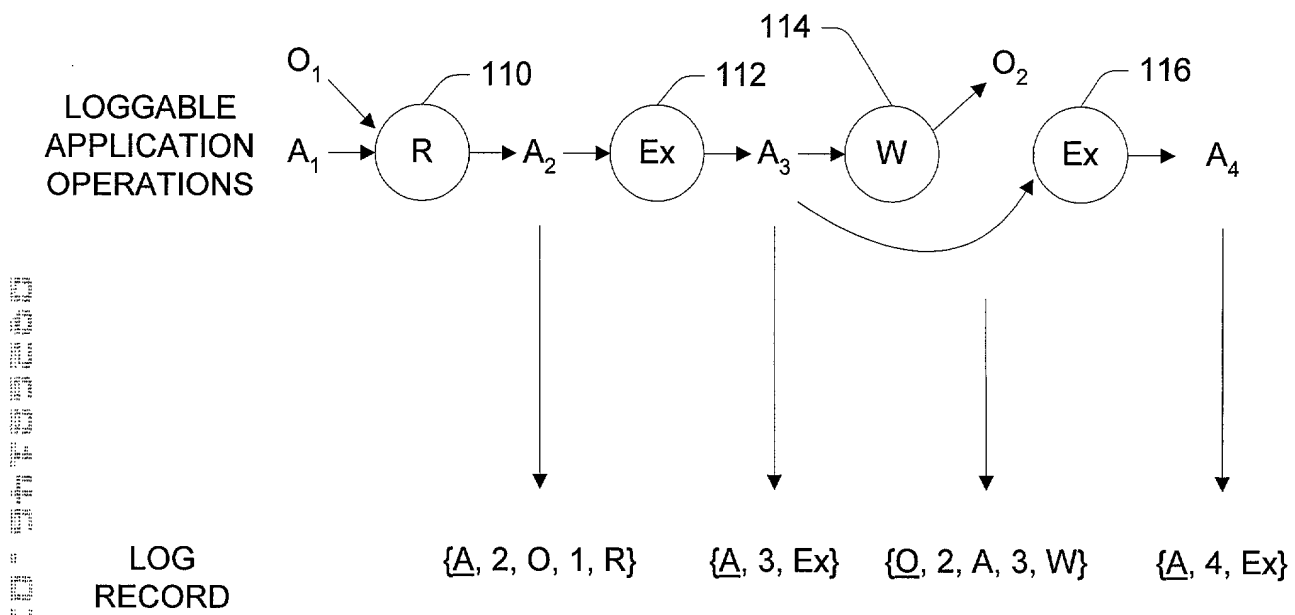
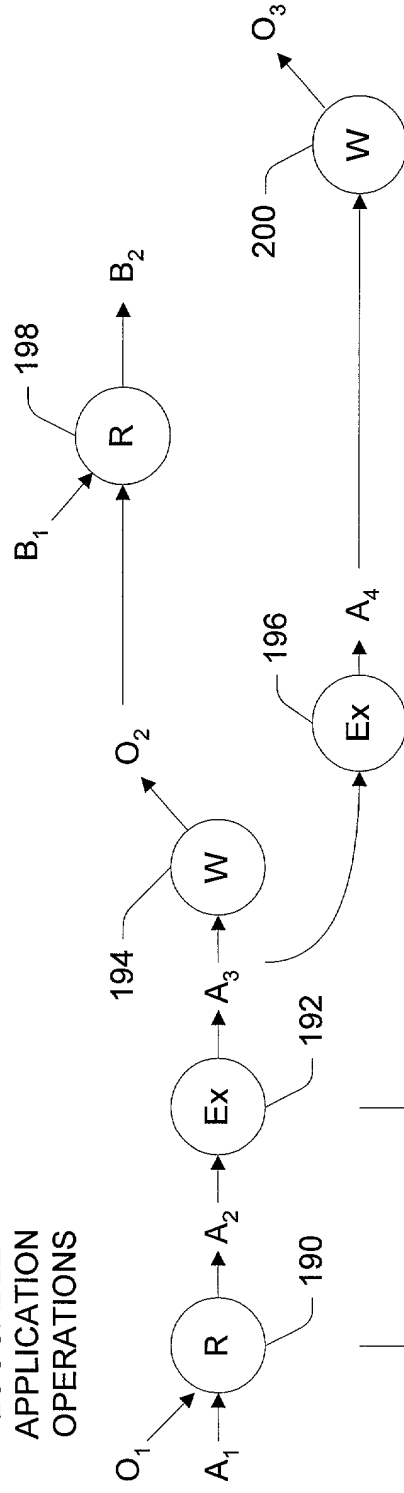


Fig. 14

LOGGABLE APPLICATION OPERATIONS



WRITE GRAPHS

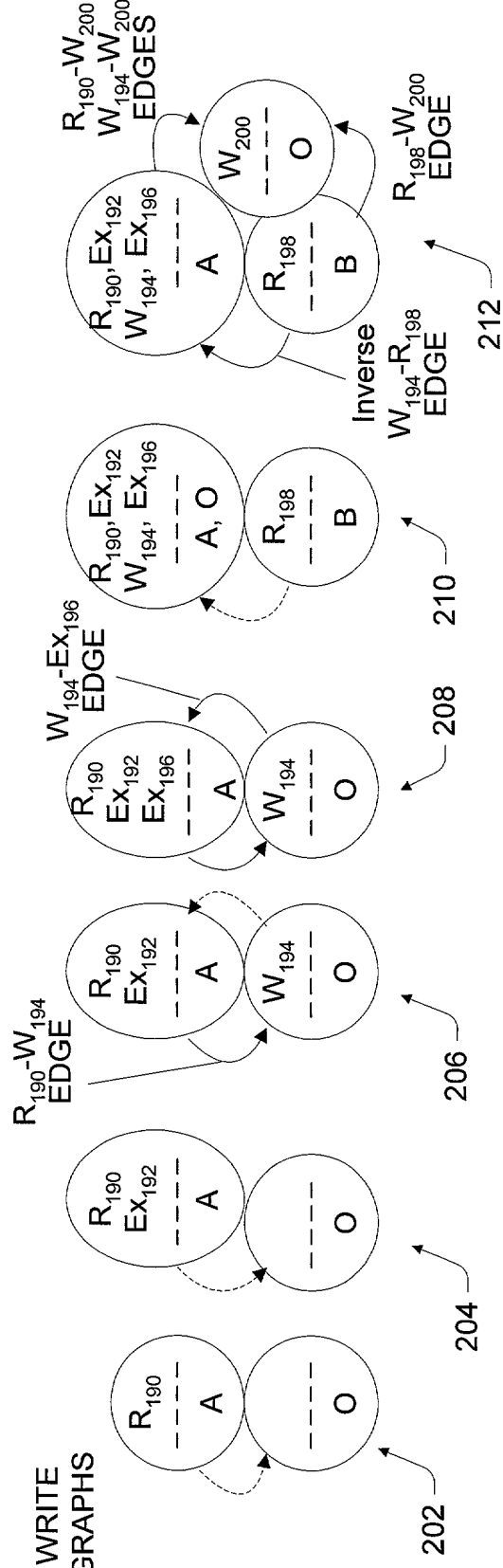


Fig. 15

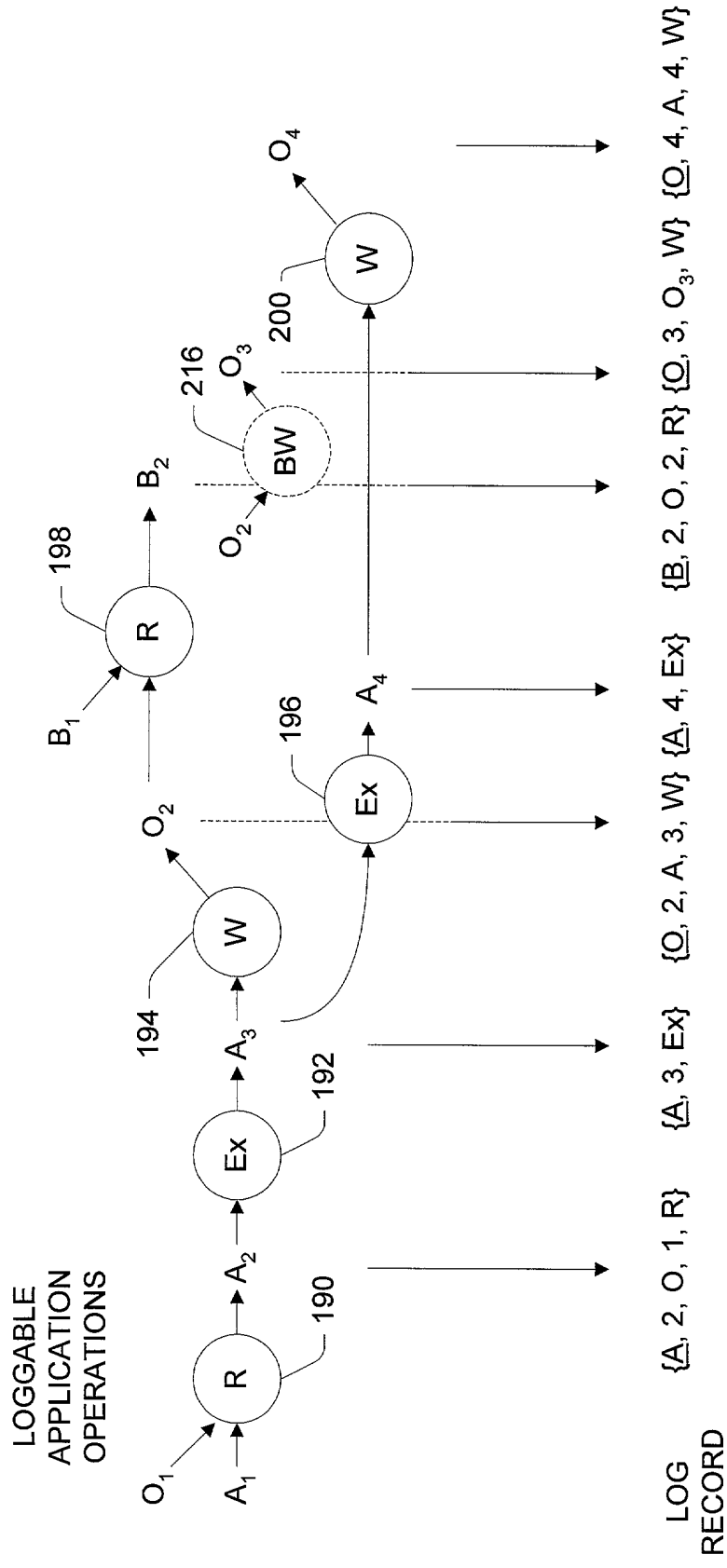


Fig. 16

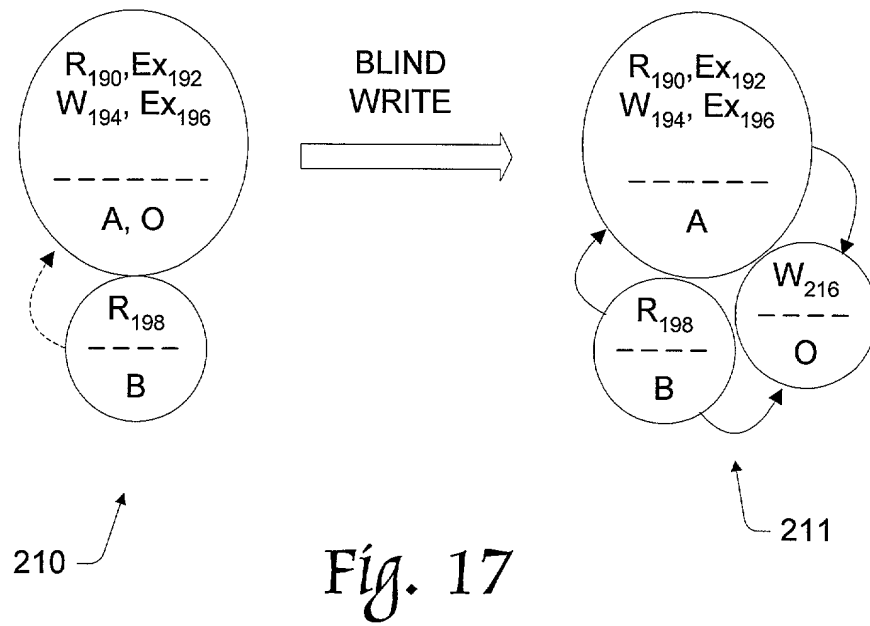


Fig. 17

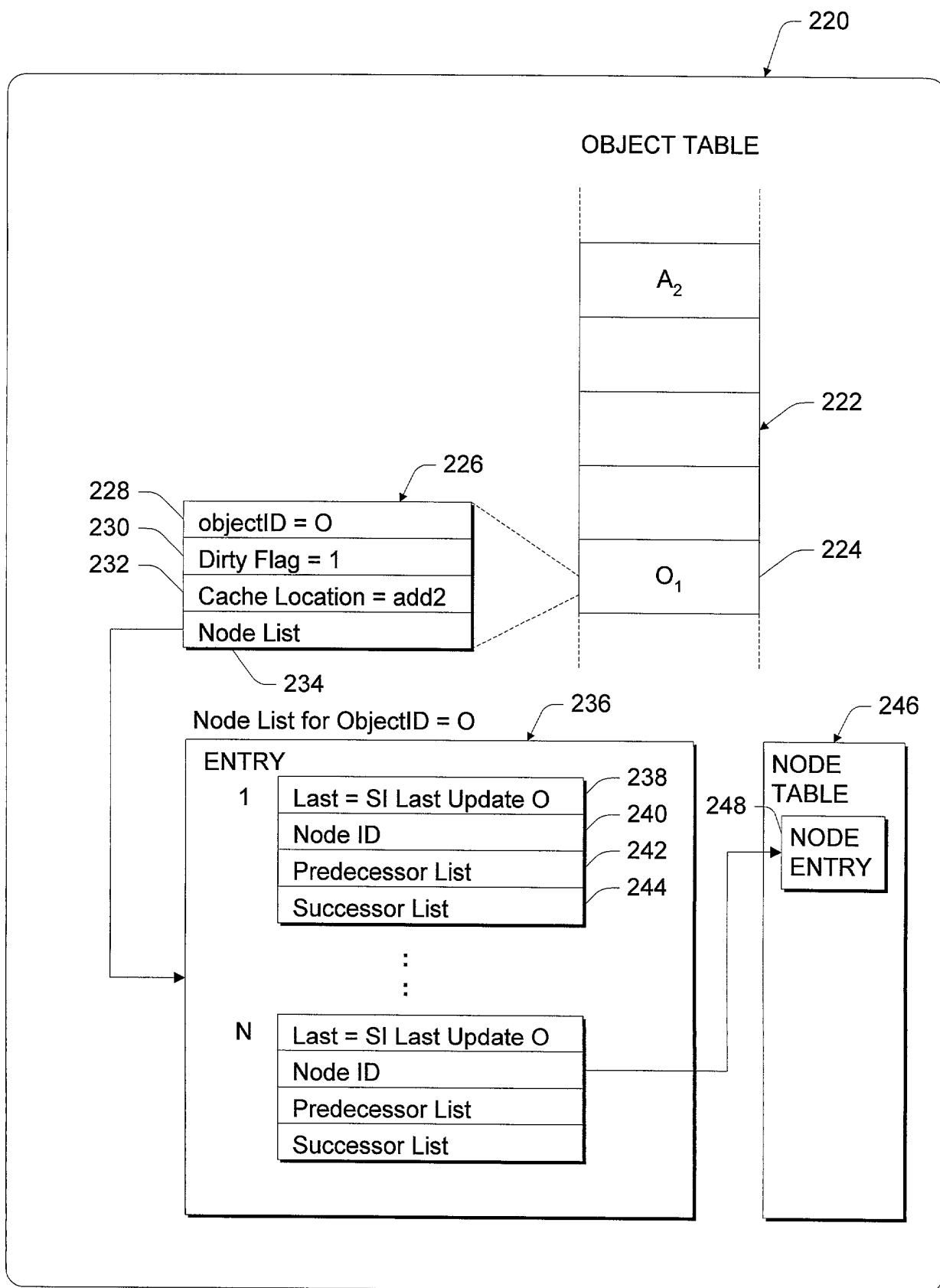


Fig. 18

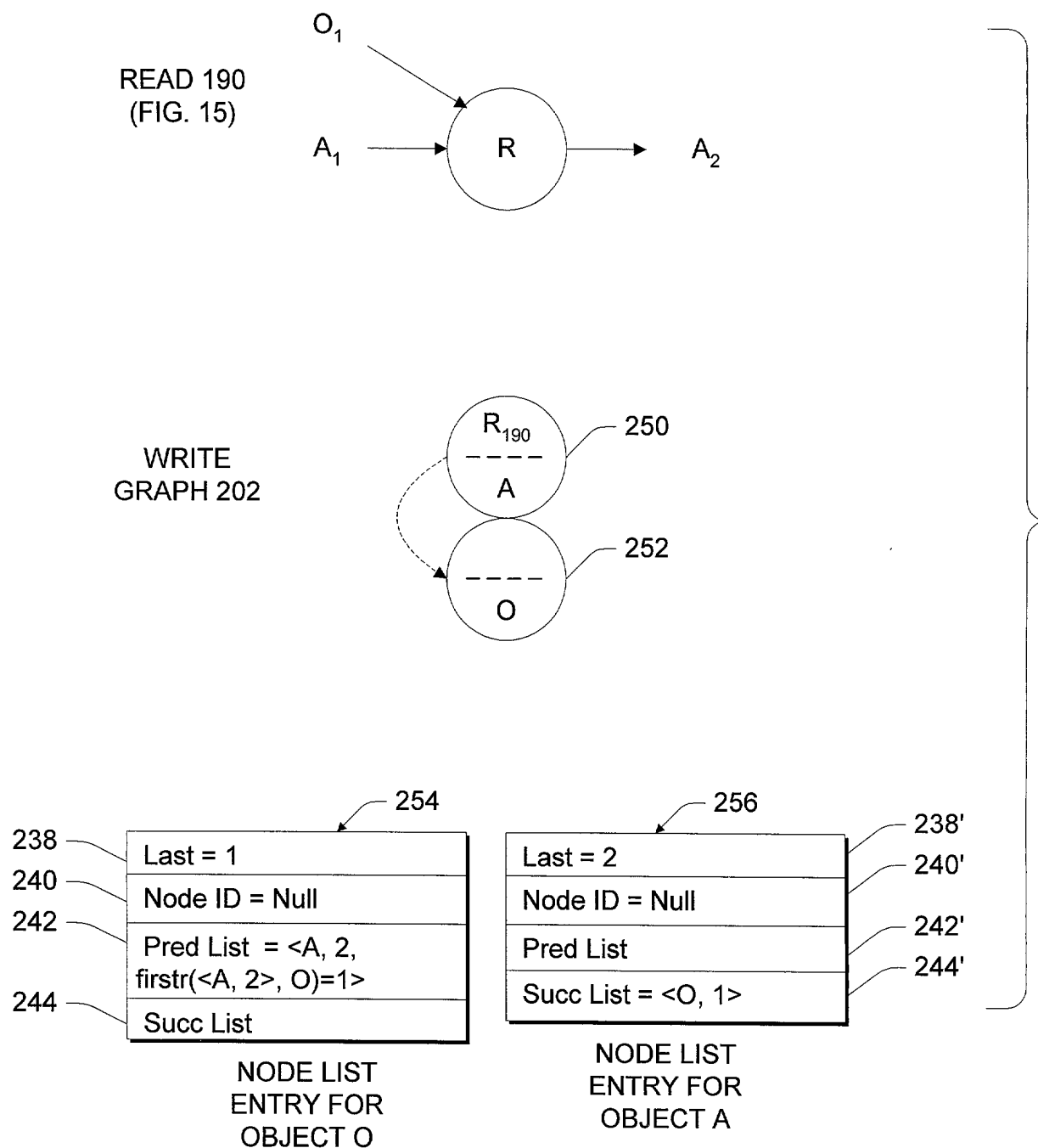
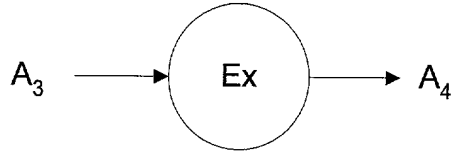
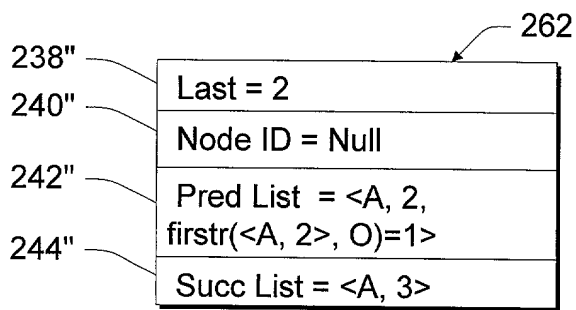
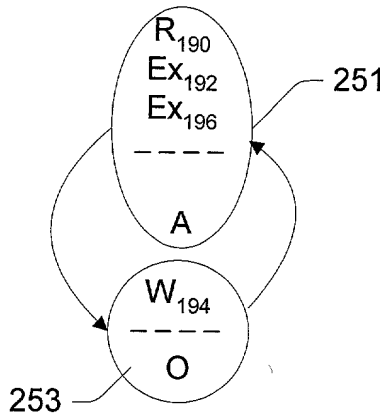


Fig. 19

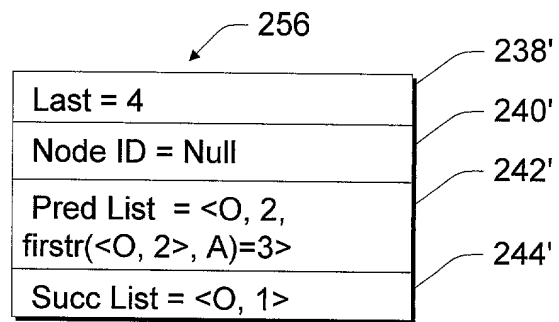
Execute 196
(FIG. 15)



WRITE
GRAPH 208



NODE LIST
ENTRY FOR
OBJECT O



NODE LIST
ENTRY FOR
OBJECT A

Fig. 20

WRITE
GRAPH 209

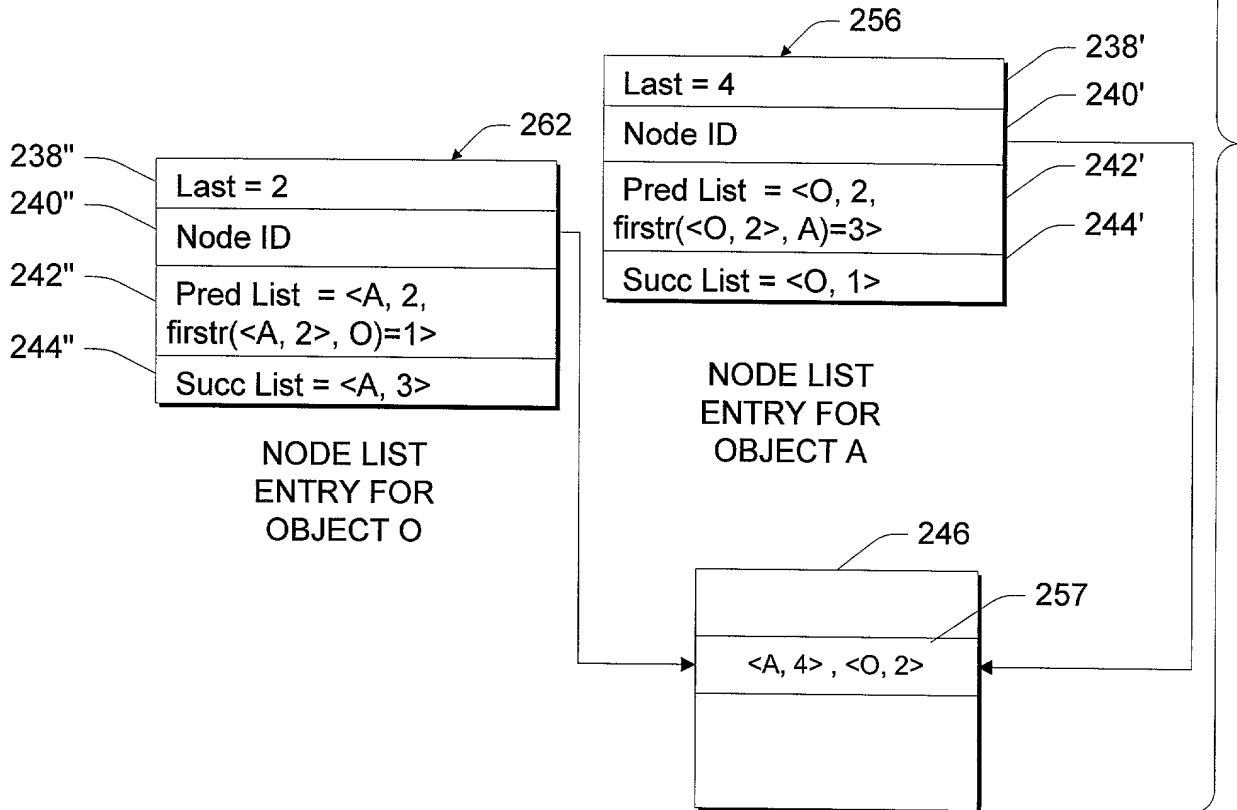
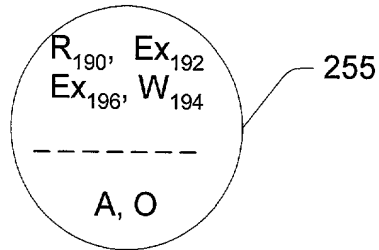
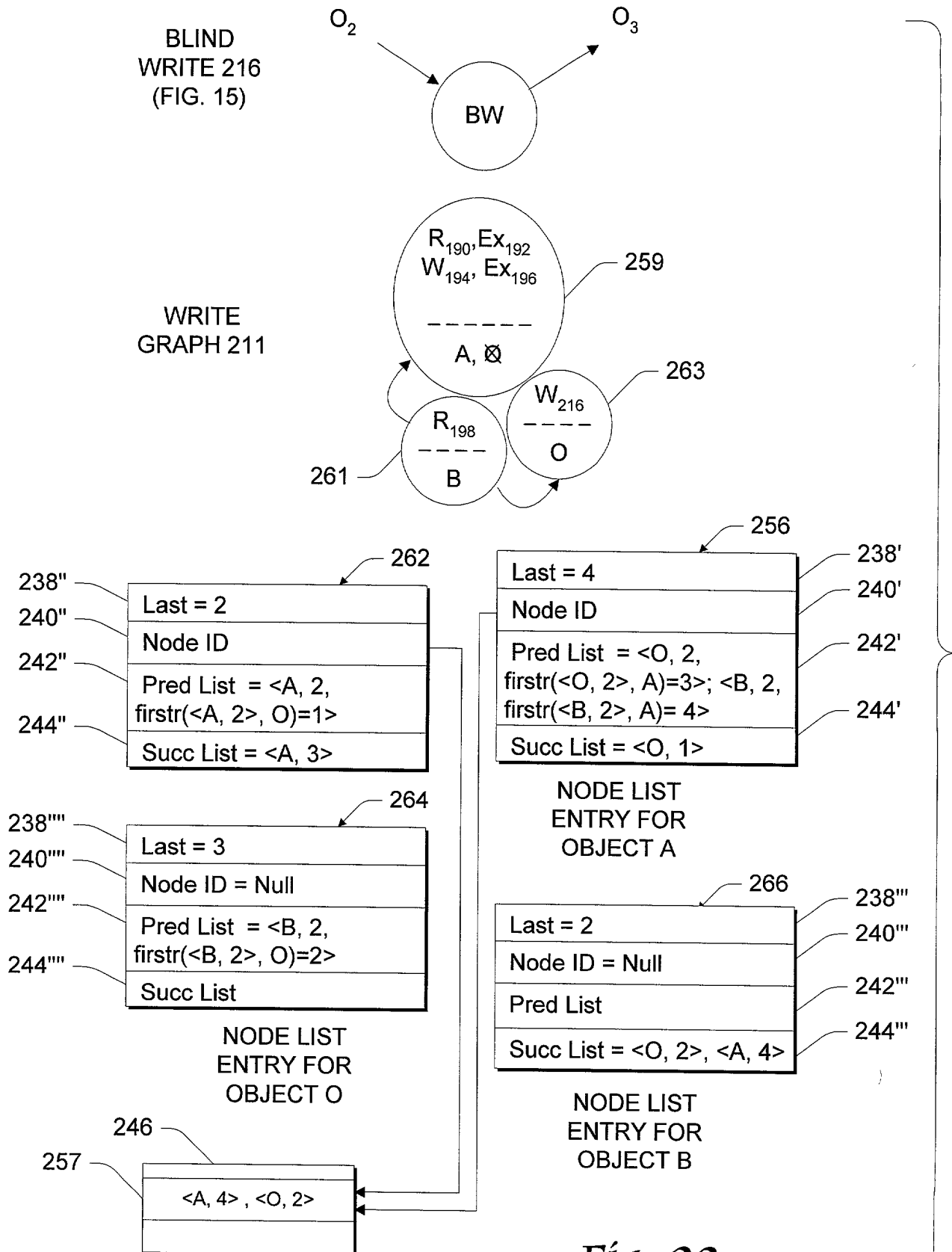
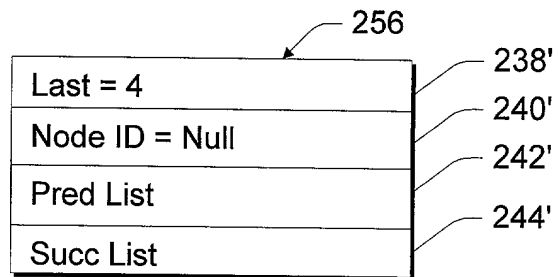
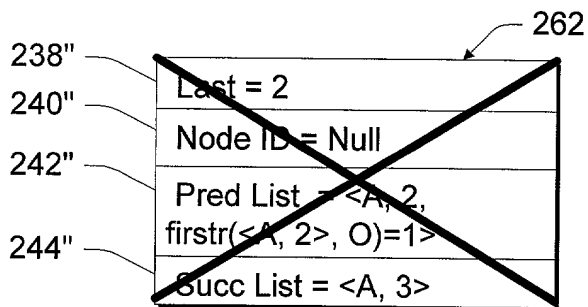
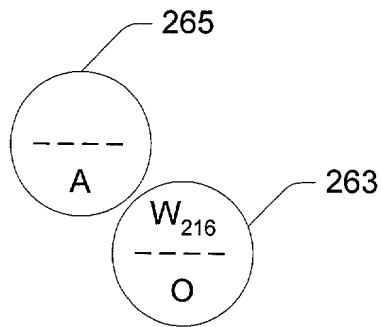


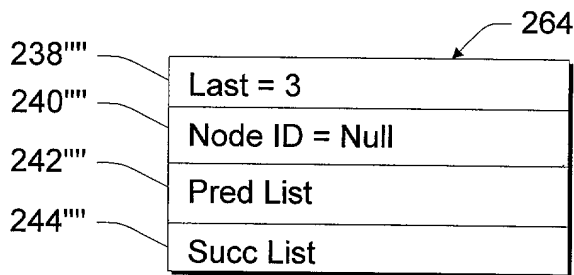
Fig. 21



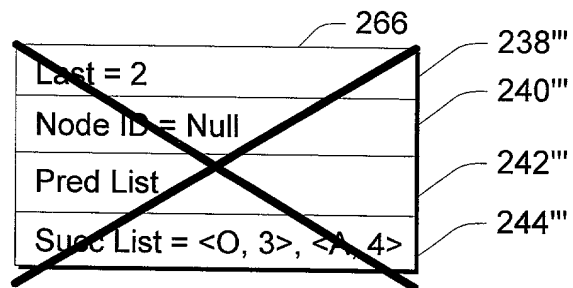
WRITE
GRAPH



NODE LIST
ENTRY FOR
OBJECT A

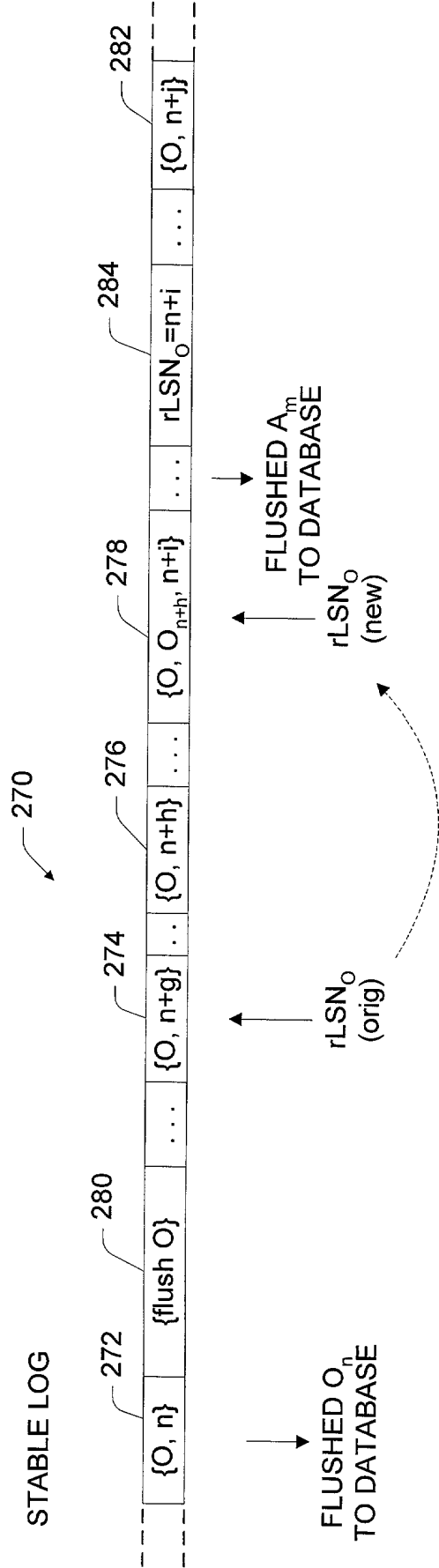
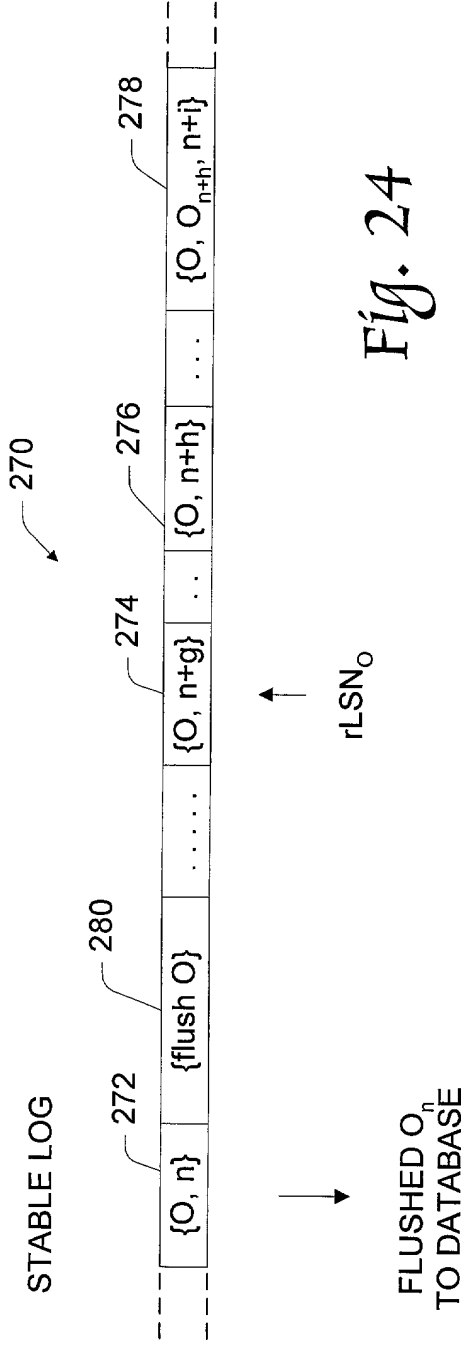


NODE LIST
ENTRY FOR
OBJECT O



NODE LIST
ENTRY FOR
OBJECT B

Fig. 23



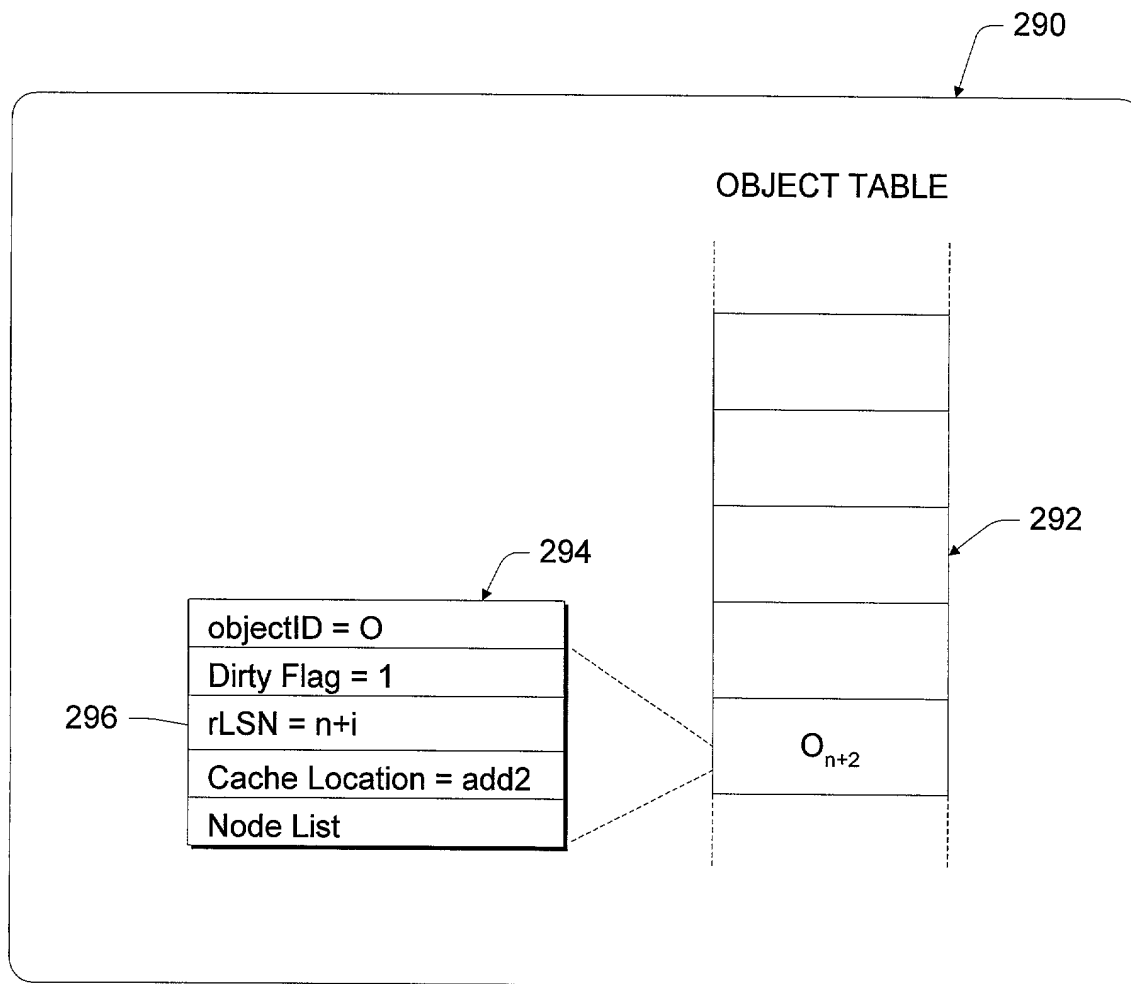


Fig. 26

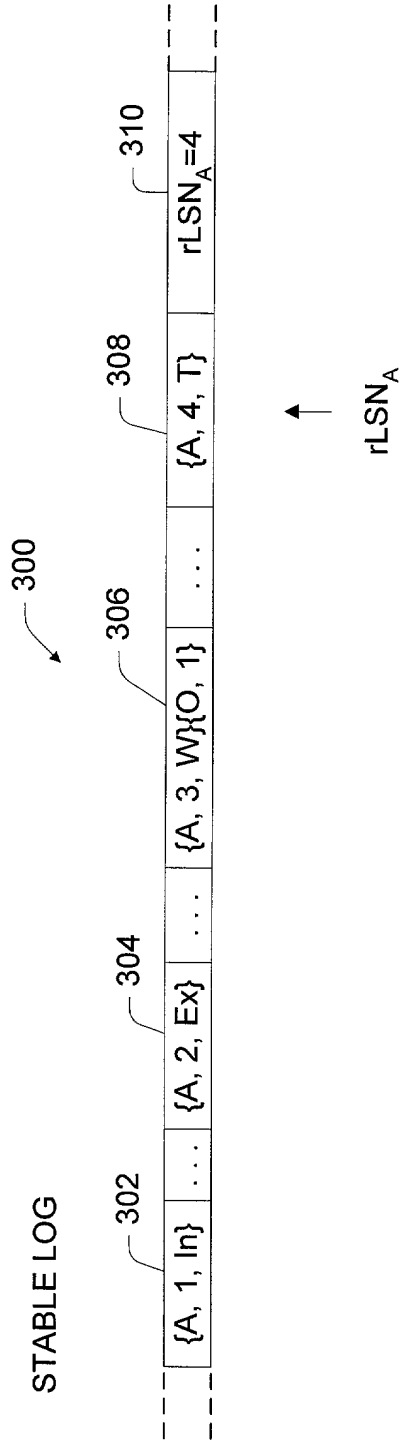


Fig. 27

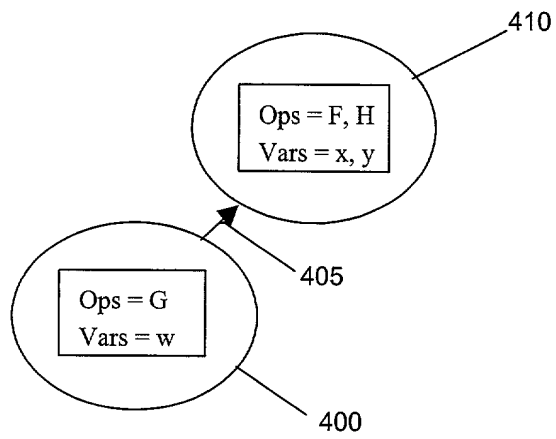


Fig. 28A

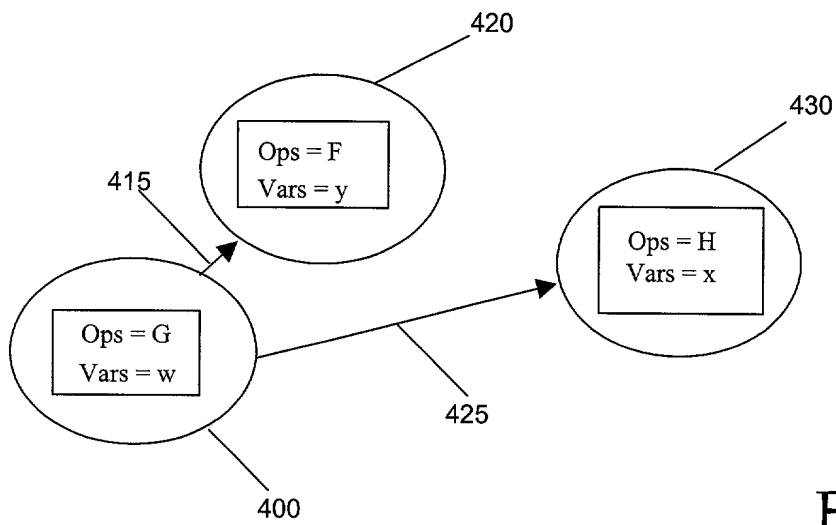


Fig. 28B

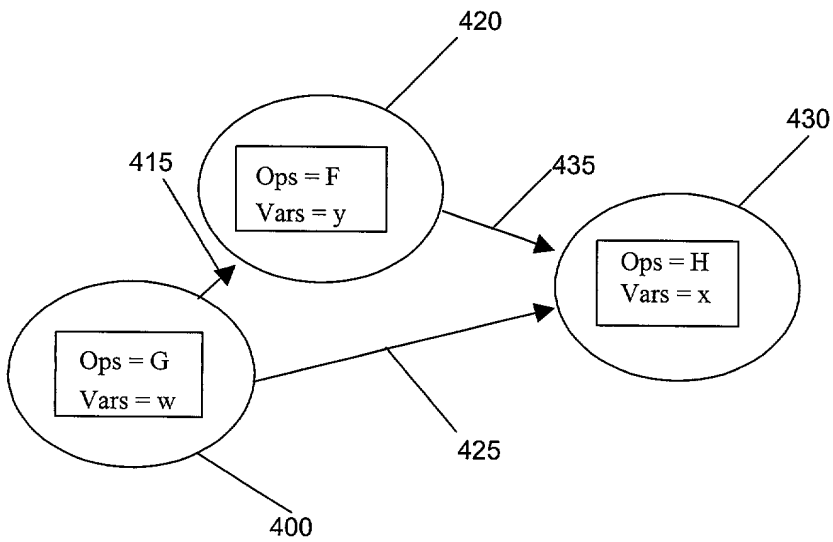


Fig. 28C

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In Re Application of:

David B. Lomet; Mark R. Tuttle

Group Art Unit: Not yet assigned

Examiner: Not yet assigned

For: DATABASE COMPUTER SYSTEM
USING LOGICAL LOGGING TO
EXTEND RECOVERY

DECLARATION AND POWER OF ATTORNEY

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name; and

I believe that I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a

☒ Utility Patent ☐ Design Patent

is sought on the invention, whose title appears above, the specification of which:

- ☒ is attached hereto.
- ☐ was filed on _____ as Serial No. _____.
- ☐ said application having been amended on _____.

I hereby state that I have reviewed and understand the contents of the above-identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose to the U.S. Patent and Trademark Office all information known to be material to the patentability of this application in accordance with 37 CFR § 1.56.

I hereby claim foreign priority benefits under 35 U.S.C. § 119(a-d) of any **foreign application(s)** for patent or inventor's certificate listed below and have also identified below any foreign application for patent or inventor's certificate having a filing date before that of any application on which priority is claimed:

Priority Claimed (If X'd)	Country	Serial Number	Date Filed
<input type="checkbox"/>	_____	_____	_____
<input type="checkbox"/>	_____	_____	_____
<input type="checkbox"/>	_____	_____	_____
<input type="checkbox"/>	_____	_____	_____

I hereby claim the benefit under 35 U.S.C. § 120 of any United States application(s) listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States application in the manner provided by the first paragraph of 35 U.S.C. § 112, I acknowledge the duty to disclose to the U.S. Patent and Trademark Office all information known to be material to patentability as defined in 37 CFR § 1.56 which became available between the filing date of the prior application and the national or PCT international filing date of this application:

Serial Number	Date Filed	Patented/Pending/Abandoned
<u>08/832,870</u>	<u>April 4, 1997</u>	<u>Pending</u>
<u>08/814,808</u>	<u>March 10, 1997</u>	<u>Pending</u>
<u>08/813,982</u>	<u>March 10, 1997</u>	<u>Patented</u>
_____	_____	_____

I hereby claim the benefit under 35 U.S.C. § 119(e) of any United States provisional application(s) listed below:

Serial Number	Date Filed
_____	_____
_____	_____

I hereby appoint the following persons of the firm of **WOODCOCK WASHBURN KURTZ MACKIEWICZ & NORRIS LLP**, One Liberty Place - 46th Floor, Philadelphia, Pennsylvania 19103 as attorney(s) and/or agent(s) to prosecute this application and to transact all business in the Patent and Trademark Office connected therewith:

	<u>Jonathan M. Waldman</u>	Reg. No. <u>40,861</u>
And	<u>Katie E. Sako</u>	Reg. No. <u>32,628</u>
	<u>Daniel D. Crouse</u>	Reg. No. <u>32,022</u>

of Microsoft Corporation, One Microsoft Way, Redmond WA 98052

Address all telephone calls and correspondence to the first-listed attorney of record at:

**WOODCOCK WASHBURN KURTZ
MACKIEWICZ & NORRIS LLP**
One Liberty Place - 46th Floor
Philadelphia PA 19103
Telephone No.: **(215) 568-3100**
Facsimile No.: **(215) 568-3439**

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

Name: David B. Lomet	
Mailing Address:	Signature
City/State of Actual Residence:	Date of Signature: _____
	Citizenship: _____

Name: Mark R. Tuttle	
Mailing Address:	Signature
City/State of Actual Residence:	Date of Signature: _____
	Citizenship: _____

Name:	
Mailing Address:	Signature
City/State of Actual Residence:	Date of Signature: _____
	Citizenship: _____